



Bayesian Deep Learning (MLSS 2019)

Yarin Gal

University of Oxford
yarin@cs.ox.ac.uk

Previously..

- ▶ Bayesian probabilistic modelling of functions
- ▶ Analytical inference of W (mean)

Today:

- ▶ Uncertainty over functions (and decomposing uncertainty)
- ▶ Scaling ideas up (approximate inference)
- ▶ Scaling up even more (stochastic approximate inference)
- ▶ Uncertainty in shallow classification models
- ▶ Stochastic approximate inference in deep NN
- ▶ Inference in very large deep models
- ▶ Real-world applications of model uncertainty

All resources (including these slides): bdl101.ml



SLIDES

Slide decks from the talks.

[SLIDE DECK 1](#)

[SLIDE DECK 2](#)



DEMO

Demoes mentioned in the slides

[UNCERTAINTY PLAYGROUND](#)

[UNCERTAINTY VISUALISATION](#)



RECAP

A quick recap of useful stuff.

[GAUSSIANS RECAP](#)



NOTATION

Notation used in the slides:.



MORE STUFF

[OATML](#)

Uncertainty over Functions

Model

- ▶ prior

$$p(w_{k,d}) = \mathcal{N}(w_{k,d}; 0, s^2); \quad W \in \mathbb{R}^{K \times D}$$

- ▶ likelihood

$$p(\mathbf{Y}|\mathbf{X}, W) = \prod_n \mathcal{N}(y_n; f^W(x_n), \sigma^2); \quad f^W(x) = W^T \phi(x)$$

- ▶ with $\phi(x)$ a K dim feature vector

Posterior

$$\begin{aligned}
 p(W|\mathbf{X}, \mathbf{Y}) &= \mathcal{N}(W; \mu', \Sigma') \\
 \Sigma' &= (\sigma^{-2} \Phi(\mathbf{X})^T \Phi(\mathbf{X}) + s^{-2} I_K)^{-1} \\
 \mu' &= \Sigma' \sigma^{-2} \Phi(\mathbf{X})^T \mathbf{Y}
 \end{aligned}$$

Predictive

$$p(y^*|x^*, \mathbf{X}, \mathbf{Y}) = \mathcal{N}(y^*; \mu'^T \phi(x^*), ?)$$

Model

- ▶ prior

$$p(w_{k,d}) = \mathcal{N}(w_{k,d}; 0, s^2); \quad W \in \mathbb{R}^{K \times D}$$

- ▶ likelihood

$$p(\mathbf{Y}|\mathbf{X}, W) = \prod_n \mathcal{N}(y_n; f^W(x_n), \sigma^2); \quad f^W(x) = W^T \phi(x)$$

- ▶ with $\phi(x)$ a K dim feature vector

Posterior

$$\begin{aligned}
 p(W|\mathbf{X}, \mathbf{Y}) &= \mathcal{N}(W; \mu', \Sigma') \\
 \Sigma' &= (\sigma^{-2} \Phi(\mathbf{X})^T \Phi(\mathbf{X}) + s^{-2} I_K)^{-1} \\
 \mu' &= \Sigma' \sigma^{-2} \Phi(\mathbf{X})^T \mathbf{Y}
 \end{aligned}$$

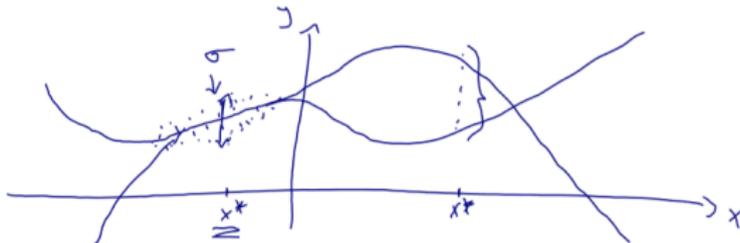
Predictive

$$p(y^*|x^*, \mathbf{X}, \mathbf{Y}) = \mathcal{N}(y^*; \mu'^T \phi(x^*), \sigma^2 + \phi(x^*)^T \Sigma' \phi(x^*))$$

$$p(y^*|x^*, X, Y) = \mathcal{N}(y^*; \mu'^T \phi(x^*), \sigma^2 + \phi(x^*)^T \Sigma' \phi(x^*))$$

Uncertainty has two components:

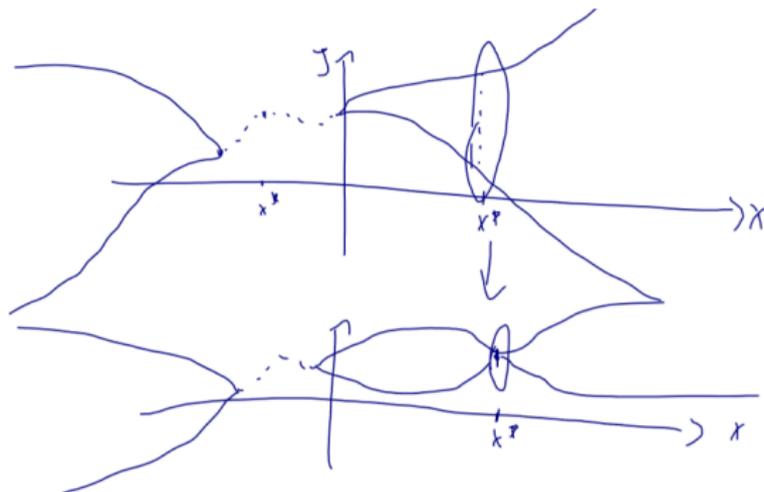
- ▶ σ^2 – from **likelihood**
- ▶ $\phi(x^*)^T \Sigma' \phi(x^*)$ – from **posterior**



$$p(y^*|x^*, X, Y) = \mathcal{N}(y^*; \mu'^T \phi(x^*), \sigma^2 + \phi(x^*)^T \Sigma' \phi(x^*))$$

Uncertainty has two components:

- ▶ σ^2 – from **likelihood**
- ▶ $\phi(x^*)^T \Sigma' \phi(x^*)$ – from **posterior**





- ▶ first term in predictive uncertainty
 $\sigma^2 + \phi(\mathbf{x}^*)^T \Sigma' \phi(\mathbf{x}^*)$
- ▶ same as likelihood σ^2 – obs noise /
corrupting additive noise eg measurement
error
- ▶ can be found via MLE rather than
assume known in advance (we'll see later)
- ▶ from Latin aleator 'dice player', from alea
'die'
 - ▶ roll a pair of dice again and again – will
not reduce uncertainty

- ▶ second term in predictive uncertainty $\sigma^2 + \phi(x^*)^T \Sigma' \phi(x^*)$
- ▶ uncertainty over **function values** before noise corruption

$$f^* = W^T \phi(x^*)$$

$$\text{Var}_{p(f^*|x^*, X, Y)}[f^*] = \phi(x^*)^T \Sigma' \phi(x^*)$$

- ▶ high for x^* “far away” from the data, even in noiseless case (ie likelihood noise is zero)
- ▶ will diminish given label for x^*
- ▶ from Ancient Greek episteme ‘knowledge, understanding’

Approximate Inference

- ▶ to evaluate predictive need to invert post cov matrix – a K by K matrix
 - ▶ difficult when K is large...
- ▶ instead, let's try to approximate posterior w a simpler dist to allow easier computations
- ▶ in approx inference we approx posterior $p(W|X, Y)$ w a different dist $q_{\theta}(W)$ param by theta
 - ▶ q also called “variational distribution”
 - ▶ θ also called “variational params”
 - ▶ technique is also known as “variational inference (VI)”
- ▶ eg q Gaussian w params $\theta = \{\mu_{VI}, \Sigma_{VI}\}$
 - ▶ $q_{\theta}(W) = \mathcal{N}(W; \mu_{VI}, \Sigma_{VI})$
 - ▶ often omit θ from subscript to avoid clutter, write $q(W)$ or q
 - ▶ often swap θ for μ, Σ back and forth

- ▶ to evaluate predictive need to invert post cov matrix – a K by K matrix
 - ▶ difficult when K is large...
- ▶ instead, let's try to approximate **posterior** w a simpler dist to allow easier computations
- ▶ in approx inference we approx **posterior** $p(W|X, Y)$ w a different dist $q_{\theta}(W)$ param by theta
 - ▶ q also called “variational distribution”
 - ▶ θ also called “variational params”
 - ▶ technique is also known as “variational inference (VI)”
- ▶ eg q Gaussian w params $\theta = \{\mu_{VI}, \Sigma_{VI}\}$
 - ▶ $q_{\theta}(W) = \mathcal{N}(W; \mu_{VI}, \Sigma_{VI})$
 - ▶ often omit θ from subscript to avoid clutter, write $q(W)$ or q
 - ▶ often swap θ for μ, Σ back and forth

- ▶ to evaluate predictive need to invert post cov matrix – a K by K matrix
 - ▶ difficult when K is large...
- ▶ instead, let's try to approximate **posterior** w a simpler dist to allow easier computations
- ▶ in approx inference we approx **posterior** $p(W|X, Y)$ w a different dist $q_{\theta}(W)$ param by theta
 - ▶ q also called “variational distribution”
 - ▶ θ also called “variational params”
 - ▶ technique is also known as “variational inference (VI)”
- ▶ eg q Gaussian w params $\theta = \{\mu_{VI}, \Sigma_{VI}\}$
 - ▶ $q_{\theta}(W) = \mathcal{N}(W; \mu_{VI}, \Sigma_{VI})$
 - ▶ often omit θ from subscript to avoid clutter, write $q(W)$ or q
 - ▶ often swap θ for μ, Σ back and forth

- ▶ to evaluate predictive need to invert post cov matrix – a K by K matrix
 - ▶ difficult when K is large...
- ▶ instead, let's try to approximate **posterior** w a simpler dist to allow easier computations
- ▶ in approx inference we approx **posterior** $p(W|X, Y)$ w a different dist $q_\theta(W)$ param by theta
 - ▶ q also called “variational distribution”
 - ▶ θ also called “variational params”
 - ▶ technique is also known as “variational inference (VI)”
- ▶ eg q Gaussian w params $\theta = \{\mu_{VI}, \Sigma_{VI}\}$
 - ▶ $q_\theta(W) = \mathcal{N}(W; \mu_{VI}, \Sigma_{VI})$
 - ▶ often omit θ from subscript to avoid clutter, write $q(W)$ or q
 - ▶ often swap θ for μ, Σ back and forth

- ▶ eg: I have posterior $p(W|X, Y) = \mathcal{N}(0, 1)$; I give you 2 approx dists

$$q_1(W) = \mathcal{N}(1, 1), \quad q_2(W) = \mathcal{N}(10, 1)$$

- ▶ which would you choose?
 - ▶ the one that gives best preds?
 - ▶ will fail: best preds are at $\mu = \mu_{\text{MLE}}, \Sigma = 0$
- ▶ need some measure of how “similar” dists are to posterior...
 - ▶ choose a measure of “similarity” between dists \tilde{D} (not necessarily a distance!)
 - ▶ then min whatever measure we commit to
 - ▶ ie if $\tilde{D}(q_1, \text{posterior}) < \tilde{D}(q_2, \text{posterior})$ then the core principle of VI says that q_1 **should be chosen over** q_2

- ▶ eg: I have posterior $p(W|X, Y) = \mathcal{N}(0, 1)$; I give you 2 approx dists

$$q_1(W) = \mathcal{N}(0, 2), \quad q_2(W) = \mathcal{N}(0, 10)$$

- ▶ which would you choose? and now?
 - ▶ the one that gives best preds?
 - ▶ will fail: best preds are at $\mu = \mu_{\text{MLE}}, \Sigma = 0$
- ▶ need some measure of how “similar” dists are to posterior...
 - ▶ choose a measure of “similarity” between dists \tilde{D} (not necessarily a distance!)
 - ▶ then min whatever measure we commit to
 - ▶ ie if $\tilde{D}(q_1, \text{posterior}) < \tilde{D}(q_2, \text{posterior})$ then the core principle of VI says that q_1 **should be chosen over** q_2

- ▶ eg: I have posterior $p(W|X, Y) = \mathcal{N}(0, 1)$; I give you 2 approx dists

$$q_1(W) = \mathcal{N}(10, 1), \quad q_2(W) = \mathcal{N}(0, 10)$$

- ▶ which would you choose? and now? ... and now?
 - ▶ the one that gives best preds?
 - ▶ will fail: best preds are at $\mu = \mu_{\text{MLE}}, \Sigma = 0$
- ▶ need some measure of how “similar” dists are to posterior...
 - ▶ choose a measure of “similarity” between dists \tilde{D} (not necessarily a distance!)
 - ▶ then min whatever measure we commit to
 - ▶ ie if $\tilde{D}(q_1, \text{posterior}) < \tilde{D}(q_2, \text{posterior})$ then the core principle of VI says that q_1 **should be chosen over** q_2

- ▶ eg: I have posterior $p(W|X, Y) = \mathcal{N}(0, 1)$; I give you 2 approx dists

$$q_1(W) = \mathcal{N}(10, 1), \quad q_2(W) = \mathcal{N}(0, 10)$$

- ▶ which would you choose? and now? ... and now?
 - ▶ the one that gives best preds?
 - ▶ will fail: best preds are at $\mu = \mu_{\text{MLE}}, \Sigma = 0$
- ▶ need some measure of how “similar” dists are to posterior...
 - ▶ choose a measure of “similarity” between dists \tilde{D} (not necessarily a distance!)
 - ▶ then min whatever measure we commit to
 - ▶ ie if $\tilde{D}(q_1, \text{posterior}) < \tilde{D}(q_2, \text{posterior})$ then the core principle of VI says that q_1 **should be chosen over** q_2

- ▶ $\tilde{D}(q_1, \text{posterior}) < \tilde{D}(q_2, \text{posterior}) \rightarrow q_1$ **should be chosen over** q_2
 - ▶ what if we have two divergences \tilde{D}_1 and \tilde{D}_2 , one saying to select q_1 and the other q_2 ?
 - ! a difference to full Bayesian inference... (where there's only one way of doing things 'correctly')
 - ▶ “from **dogmatic** Bayes to **pragmatic** Bayes”;
 - ▶ often choose \tilde{D} that is mathematically convenient
- ▶ eg Kullback Leibler

$$\text{KL}(q, p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

- ▶ $\tilde{D}(q_1, \text{posterior}) < \tilde{D}(q_2, \text{posterior}) \rightarrow q_1$ **should be chosen over** q_2
 - ▶ what if we have two divergences \tilde{D}_1 and \tilde{D}_2 , one saying to select q_1 and the other q_2 ?
 - ! a difference to full Bayesian inference... (where there's only one way of doing things 'correctly')
 - ▶ “from **dogmatic** Bayes to **pragmatic** Bayes”;
 - ▶ often choose \tilde{D} that is mathematically convenient
- ▶ eg Kullback Leibler

$$\text{KL}(q, p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

- ▶ K dim discrete prob vectors q, p : $KL(q, p) = \sum_k q_k \log q_k/p_k$
- ▶ when the two dists are the **same** we get **exactly 0**
- ▶ when the two dists are **different** the divergence is **positive**
- ▶ KL is **not symmetric**
- ▶ if q_k is zero it is **ignored** in KL
- ▶ whenever $q_k > 0$ it must be that $p_k > 0$ for the KL to be **finite**
- ▶ Homework: find examples for all properties; eg
 $q = [1/8, 3/8, 4/8], p = [3/8, 4/8, 1/8]$;

- ▶ K dim discrete prob vectors q, p : $KL(q, p) = \sum_k q_k \log q_k/p_k$
- ▶ when the two dists are the **same** we get **exactly 0**
- ▶ when the two dists are **different** the divergence is **positive**
- ▶ KL is **not symmetric**
- ▶ if q_k is zero it is **ignored** in KL
- ▶ whenever $q_k > 0$ it must be that $p_k > 0$ for the KL to be **finite**
- ▶ Homework: find examples for all properties; eg
 $q = [1/8, 3/8, 4/8], p = [3/8, 4/8, 1/8]$;

- ▶ K dim discrete prob vectors q, p : $KL(q, p) = \sum_k q_k \log q_k/p_k$
- ▶ when the two dists are the **same** we get **exactly 0**
- ▶ when the two dists are **different** the divergence is **positive**
- ▶ KL is **not symmetric**
- ▶ if q_k is zero it is **ignored** in KL
- ▶ whenever $q_k > 0$ it must be that $p_k > 0$ for the KL to be **finite**
- ▶ Homework: find examples for all properties; eg
 $q = [1/8, 3/8, 4/8], p = [3/8, 4/8, 1/8];$

- ▶ K dim discrete prob vectors q, p : $KL(q, p) = \sum_k q_k \log q_k/p_k$
- ▶ when the two dists are the **same** we get **exactly 0**
- ▶ when the two dists are **different** the divergence is **positive**
- ▶ KL is **not symmetric**
- ▶ if q_k is zero it is **ignored** in KL
- ▶ whenever $q_k > 0$ it must be that $p_k > 0$ for the KL to be **finite**
- ▶ Homework: find examples for all properties; eg
 $q = [1/8, 3/8, 4/8], p = [3/8, 4/8, 1/8];$

- ▶ K dim discrete prob vectors q, p : $KL(q, p) = \sum_k q_k \log q_k/p_k$
- ▶ when the two dists are the **same** we get **exactly 0**
- ▶ when the two dists are **different** the divergence is **positive**
- ▶ KL is **not symmetric**
- ▶ if q_k is zero it is **ignored** in KL
- ▶ whenever $q_k > 0$ it must be that $p_k > 0$ for the KL to be **finite**
- ▶ Homework: find examples for all properties; eg
 $q = [1/8, 3/8, 4/8], p = [3/8, 4/8, 1/8];$

- ▶ K dim discrete prob vectors q, p : $KL(q, p) = \sum_k q_k \log q_k/p_k$
- ▶ when the two dists are the **same** we get **exactly 0**
- ▶ when the two dists are **different** the divergence is **positive**
- ▶ KL is **not symmetric**
- ▶ if q_k is zero it is **ignored** in KL
- ▶ whenever $q_k > 0$ it must be that $p_k > 0$ for the KL to be **finite**
- ▶ Homework: find examples for all properties; eg
 $q = [1/8, 3/8, 4/8], p = [3/8, 4/8, 1/8];$

- ▶ K dim discrete prob vectors q, p : $KL(q, p) = \sum_k q_k \log q_k/p_k$
- ▶ when the two dists are the **same** we get **exactly 0**
- ▶ when the two dists are **different** the divergence is **positive**
- ▶ KL is **not symmetric**
- ▶ if q_k is zero it is **ignored** in KL
- ▶ whenever $q_k > 0$ it must be that $p_k > 0$ for the KL to be **finite**
- ▶ Homework: find examples for all properties; eg
 $q = [1/8, 3/8, 4/8], p = [3/8, 4/8, 1/8];$

What if we want to approx cnts rv like W ?

- ▶ $q(x) = N(x; \mu_0, s_0^2)$, $p(x) = N(x; \mu_1, s_1^2)$; KL for Gaussians:

$$\text{KL}(q, p) = 1/2(s_1^{-2}s_0^2 + s_1^{-2}(\mu_1 - \mu_0)^2 - 1 + \log(s_1^2/s_0^2))$$

- ▶ nice property: if X_1 and X_2 are independent under p and q then
 $\text{KL}(q(X_1, X_2), p(X_1, X_2)) = \text{KL}(q(X_1), p(X_1)) + \text{KL}(q(X_2), p(X_2))$
- ▶ multivariate diagonal Gaussians (K dims):
 write $\mathbf{x} = [x_1, \dots, x_K]$

$$q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mu_0, S_0) \quad \text{with } S_0 = \text{diag}([s_{01}^2, \dots, s_{0K}^2])$$

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mu_1, S_1) \quad \text{with } S_1 = \text{diag}([s_{11}^2, \dots, s_{1K}^2])$$

Then from indep of x_1, \dots, x_K :

$$\text{KL}(q, p) = \sum_k 1/2(s_{1k}^{-2}s_{0k}^2 + s_{1k}^{-2}(\mu_{1k} - \mu_{0k})^2 - 1 + \log(s_{1k}^2/s_{0k}^2))$$

What if we want to approx cnts rv like W ?

- ▶ $q(x) = N(x; \mu_0, s_0^2)$, $p(x) = N(x; \mu_1, s_1^2)$; KL for Gaussians:

$$\text{KL}(q, p) = 1/2(s_1^{-2}s_0^2 + s_1^{-2}(\mu_1 - \mu_0)^2 - 1 + \log(s_1^2/s_0^2))$$

- ▶ nice property: if X_1 and X_2 are independent under p and q then
 $\text{KL}(q(X_1, X_2), p(X_1, X_2)) = \text{KL}(q(X_1), p(X_1)) + \text{KL}(q(X_2), p(X_2))$
- ▶ multivariate diagonal Gaussians (K dims):
 write $\mathbf{x} = [x_1, \dots, x_K]$

$$q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mu_0, S_0) \quad \text{with } S_0 = \text{diag}([s_{01}^2, \dots, s_{0K}^2])$$

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mu_1, S_1) \quad \text{with } S_1 = \text{diag}([s_{11}^2, \dots, s_{1K}^2])$$

Then from indep of x_1, \dots, x_K :

$$\text{KL}(q, p) = \sum_k 1/2(s_{1k}^{-2}s_{0k}^2 + s_{1k}^{-2}(\mu_{1k} - \mu_{0k})^2 - 1 + \log(s_{1k}^2/s_{0k}^2))$$

What if we want to approx cnts rv like W ?

- ▶ $q(x) = N(x; \mu_0, s_0^2)$, $p(x) = N(x; \mu_1, s_1^2)$; KL for Gaussians:

$$\text{KL}(q, p) = 1/2(s_1^{-2}s_0^2 + s_1^{-2}(\mu_1 - \mu_0)^2 - 1 + \log(s_1^2/s_0^2))$$

- ▶ nice property: if X_1 and X_2 are independent under p and q then

$$\text{KL}(q(X_1, X_2), p(X_1, X_2)) = \text{KL}(q(X_1), p(X_1)) + \text{KL}(q(X_2), p(X_2))$$

- ▶ multivariate diagonal Gaussians (K dims):

write $\mathbf{x} = [x_1, \dots, x_K]$

$$q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mu_0, \mathbf{S}_0) \quad \text{with } \mathbf{S}_0 = \text{diag}([s_{01}^2, \dots, s_{0K}^2])$$

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mu_1, \mathbf{S}_1) \quad \text{with } \mathbf{S}_1 = \text{diag}([s_{11}^2, \dots, s_{1K}^2])$$

Then from indep of x_1, \dots, x_K :

$$\text{KL}(q, p) = \sum_k 1/2(s_{1k}^{-2}s_{0k}^2 + s_{1k}^{-2}(\mu_{1k} - \mu_{0k})^2 - 1 + \log(s_{1k}^2/s_{0k}^2))$$

- ▶ want to approx $p(W|X, Y)$ using some $q_\theta(W)$
- ▶ min

$$\text{KL}(q_\theta(W), p(W|X, Y))$$

wrt θ (remember def $\text{KL}(q, p) = \int q(x) \log \frac{q(x)}{p(x)} dx$)

$$\begin{aligned} & \text{KL}(q(w), p(w|x, y)) \\ &= \int q(w) \log \frac{d(w)}{p(w|x, y)} dw \\ & \quad \frac{p(y|w, \kappa) p(w)}{p(y|x)} \end{aligned}$$

$$\begin{aligned} & \text{KL}(q(w), p(w|x, y)) \\ &= \int q(w) \log \frac{q(w)}{p(w|x, y)} dw \\ & \quad \frac{p(y|w, x)p(w)}{p(y|x)} \\ &= \int q(w) \log \frac{q(w)p(y|x)}{p(y|w, x)p(w)} dw \end{aligned}$$

$$KL(q(w), p(w|x, y))$$

$$= \int q(w) \log \frac{d(w)}{p(w|x, y)} dw$$

$$= \int q(w) \log \frac{p(y|w, x)p(w)}{p(y|x)} dw$$

$$= \int q(w) \log \frac{q(w)p(y|x)}{p(y|w, x)p(w)} dw$$

$$= \int q(w) \log \frac{1}{p(y|w, x)} dw + \underbrace{\int q(w) \log \frac{q(w)}{p(w)} dw}_{\geq KL(q(w), p(w))} + \underbrace{\int q(w) \log p(y|x) dw}_{\text{constant}}$$

$\underbrace{\int q(w) \log \frac{1}{p(y|w, x)} dw}_{= -\log p(\cdot)}$

$$KL(q(w), p(w|x, y))$$

$$= \int q(w) \log \frac{d(w)}{p(w|x, y)} dw$$

$$= \int q(w) \log \frac{p(y|w, x)p(w)}{p(y|x)} dw$$

$$= \int q(w) \log \frac{q(w)p(y|x)}{p(y|w, x)p(w)} dw$$

$$= \int q(w) \log \frac{1}{p(y|w, x)} dw + \underbrace{\int q(w) \log \frac{q(w)}{p(w)} dw}_{= KL(q(w), p(w))} + \underbrace{\int q(w) \log p(y|x) dw}_{= \log p(y|x)}$$

$$= - \int q(w) \log p(y|w, x) dw + KL(q(w), p(w)) + \log p(y|x)$$

KL for approx inference

$$KL(q(w), p(w|x, y))$$

$$= \int q(w) \log \frac{q(w)}{\underbrace{p(w|x, y)}} dw$$
$$\frac{p(y|w, x)p(w)}{p(y|x)}$$

$$= \int q(w) \log \frac{q(w)p(y|x)}{p(y|w, x)p(w)} dw$$

$$= \int q(w) \log \frac{1}{p(y|w, x)} dw + \underbrace{\int q(w) \log \frac{q(w)}{p(w)} dw}_{= KL(q(w), p(w))} + \underbrace{\int q(w) \log p(y|x) dw}$$

$= -\log p(\dots)$ $= KL(q(w), p(w))$ $= \int q(w) \log p(y|x) dw$

$$= - \int q(w) \log p(y|w, x) dw + KL(q(w), p(w)) + \log p(y|x)$$

$$\Rightarrow \underbrace{\int q(w) \log p(y|w, x) dw}_{\mathbb{E}_{x,p} \log \text{lik}} - \underbrace{KL(q(w), p(w))}_{KL \text{ to prior}} = \underbrace{\log p(y|x)}_{\log \text{ evidence}} - \underbrace{KL(q(w), p(w|x, y))}_{\text{ELBO}}$$

KL for approx inference

$$KL(q(w), p(w|x, y))$$

$$= \int q(w) \log \frac{q(w)}{p(w|x, y)} dw$$

$$= \int q(w) \log \frac{q(w) p(y|x)}{p(y|w, x) p(w)} dw$$

$$= \int q(w) \log \frac{q(w) p(y|x)}{p(y|w, x) p(w)} dw$$

$$= \int q(w) \log \frac{1}{p(y|w, x)} dw + \underbrace{\int q(w) \log \frac{q(w)}{p(w)} dw}_{= KL(q(w), p(w))} + \underbrace{\int q(w) \log p(y|x) dw}_{}$$

$$= - \int q(w) \log p(y|w, x) dw + KL(q(w), p(w)) + \log p(y|x)$$

$$\int_{LHS} \leq \int_{RHS}$$

$$\Rightarrow \underbrace{\int q(w) \log p(y|x, w) dw}_{\text{exp log lik}} - \underbrace{KL(q(w), p(w))}_{\text{KL to prior}} \leq \underbrace{\log p(y|x)}_{\text{log evidence}} - \underbrace{KL(q(w), p(w|x, y))}_{}$$

(ELBO)

- ▶ want to approx $p(W|X, Y)$ using some $q_\theta(W)$

- ▶ min

$$\text{KL}(q_\theta(W), p(W|X, Y))$$

wrt θ (remember def $\text{KL}(q, p) = \int q(x) \log \frac{q(x)}{p(x)} dx$)

- ▶ $\log p(Y|X) \geq \int q(W) \log p(Y|X, W) dW - \text{KL}(q(W), p(W))$

- ▶ pops out a bound on evidence for free
- ▶ also called “evidence lower bound” (**ELBO**)
- ▶ min KL to posterior = max ELBO

- ▶ what does it mean to max ELBO?

- ▶ **first term**: how well we “explain the data”; if possible, q should put all mass at MLE!
- ▶ **second term**: how close we are to the prior (get simplest q that can still explain data well); if possible, q should be prior itself!

- ▶ want to approx $p(W|X, Y)$ using some $q_\theta(W)$

- ▶ min

$$\text{KL}(q_\theta(W), p(W|X, Y))$$

wrt θ (remember def $\text{KL}(q, p) = \int q(x) \log \frac{q(x)}{p(x)} dx$)

- ▶ $\log p(Y|X) \geq \int q(W) \log p(Y|X, W) dW - \text{KL}(q(W), p(W))$
 - ▶ pops out a bound on evidence for free
 - ▶ also called “evidence lower bound” (**ELBO**)
 - ▶ min KL to posterior = max ELBO
- ▶ what does it mean to max ELBO?
 - ▶ **first term**: how well we “explain the data”; if possible, q should put all mass at MLE!
 - ▶ **second term**: how close we are to the prior (get simplest q that can still explain data well); if possible, q should be prior itself!

- ▶ max

$$\int q_{\theta}(W) \log p(Y|X, W) dW - \text{KL}(q_{\theta}(W), p(W))$$

wrt θ

- ▶ which terms can we compute?
 - ▶ for Gaussian prior and q , can compute KL to prior
 - ▶ for Gaussian lik can compute expected log lik as well (analytic – try this at home using tools from earlier!)
 - ▶ but in more complicated likelihoods (like in **classification**) can't eval above...
 - ▶ for this we'll look at **stochastic** approximate inference

- ▶ max

$$\int q_{\theta}(W) \log p(Y|X, W) dW - \text{KL}(q_{\theta}(W), p(W))$$

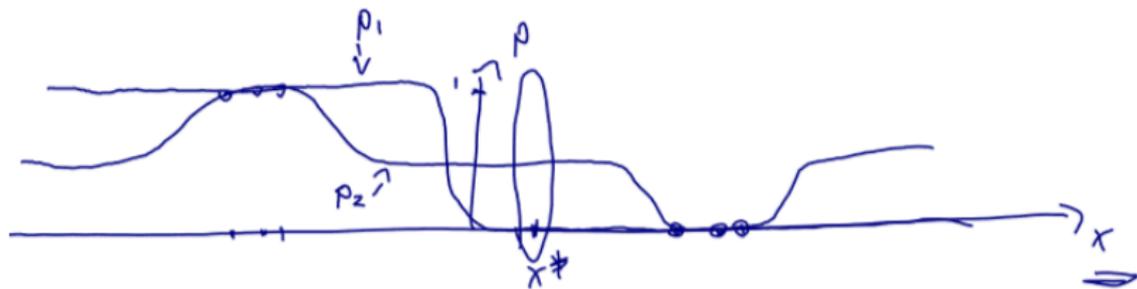
wrt θ

- ▶ which terms can we compute?
 - ▶ for Gaussian prior and q , can compute KL to prior
 - ▶ for Gaussian lik can compute expected log lik as well (analytic – try this at home using tools from earlier!)
 - ▶ but in more complicated likelihoods (like in **classification**) can't eval above...
 - ▶ for this we'll look at **stochastic** approximate inference

Stochastic Approximate Inference

Let's try to do a classification task

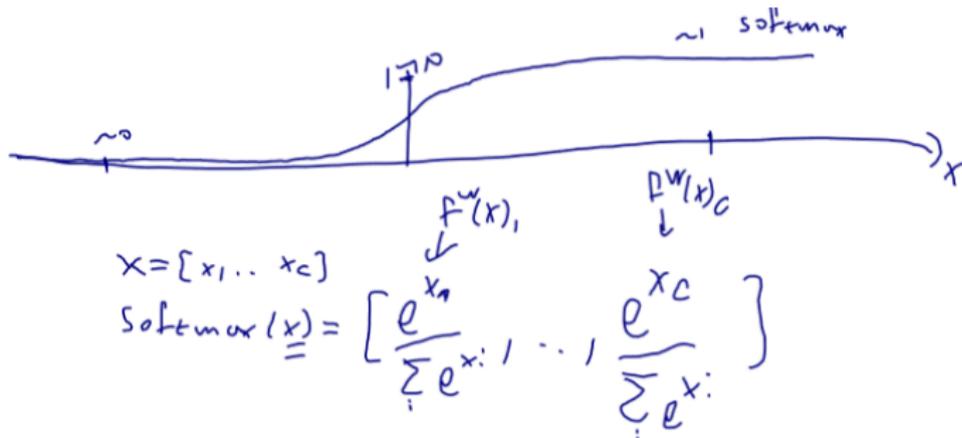
- ▶ want to get notion of epistemic uncertainty in classification



- ▶ generative story
 - ▶ Nature chose function $p(x) : \mathbb{R}^Q \rightarrow [0, 1]^C$
 - ▶ $p(x)$ a prob vector as a function of x
 - ▶ eg p softmax func
 - ▶ for $n = 1..N$ generate label $y_n \sim \text{Categorical}(p(x_n))$
- ▶ encode y_n as a one hot vector \mathbf{y}_n (eg $[0, 0, 1, 0]$ with $C = 4$ classes and $y_n = 2$)

Let's try to do a classification task

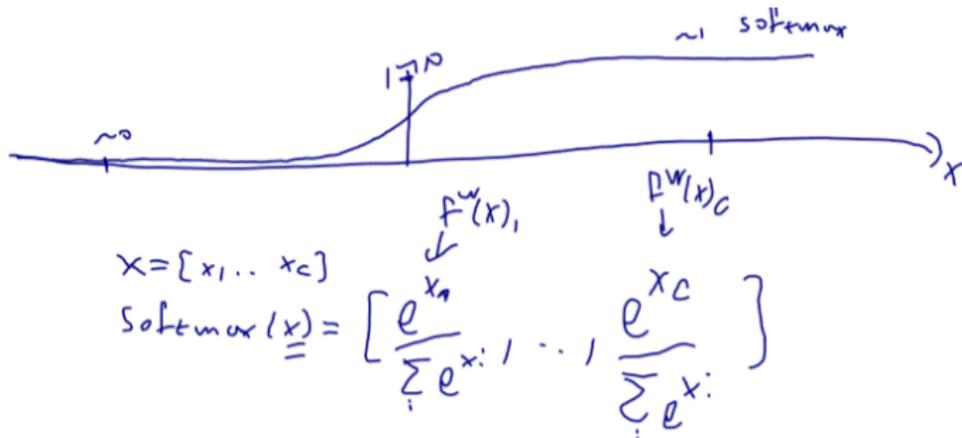
- ▶ want to get notion of epistemic uncertainty in classification
- ▶ generative story
 - ▶ Nature chose function $p(x) : \mathbb{R}^Q \rightarrow [0, 1]^C$
 - ▶ $p(x)$ a prob vector as a function of x
 - ▶ eg p softmax func
 - ▶ for $n = 1..N$ generate label $y_n \sim \text{Categorical}(p(x_n))$



- ▶ encode y_n as a one hot vector \mathbf{y}_n (eg $[0, 0, 1, 0]$ with $C = 4$ classes and $y_n = 2$)

Let's try to do a classification task

- ▶ want to get notion of epistemic uncertainty in classification
- ▶ generative story
 - ▶ Nature chose function $p(x) : \mathbb{R}^Q \rightarrow [0, 1]^C$
 - ▶ $p(x)$ a prob vector as a function of x
 - ▶ eg p softmax func
 - ▶ for $n = 1..N$ generate label $y_n \sim \text{Categorical}(p(x_n))$



- ▶ encode y_n as a one hot vector \mathbf{y}_n (eg $[0, 0, 1, 0]$ with $C = 4$ classes and $y_n = 2$)

Model:

- ▶ likelihood
 - ▶ model prob func by **function** $p^W(x)$ with W a K by C matrix; then lik is **def'd** as elem c in prob vec

$$p(y = c | x, W) = p^W(x)_c,$$

$$\begin{aligned} p(Y|X, W) &= \prod_n p^W(x_n)_{y_n=c} \\ &= \prod_n \mathbf{y}_n^T p^W(x_n) \end{aligned}$$

- ▶ prior over W
 - ▶ vectorise W (still write W instead of $\text{vec}(W)$)
 - ▶ same prior as before:

$$p(W) = \mathcal{N}(W; 0_{CK}, s^2 I_{CK})$$

For each c :

$$w_c : k \times 1$$

write

$$W = [w_1, \dots, w_C] \\ : k \times C$$

Def 'logits' =

$$f^W(x) =$$

$$[w_1^T \phi(x), \dots, w_C^T \phi(x)]$$

Def prob func

$$p^W(x) =$$

$$\text{Softmax}(f^W(x))$$

Model:

- ▶ likelihood
 - ▶ model prob func by **function** $p^W(x)$ with W a K by C matrix; then lik is **def'd** as elem c in prob vec

$$p(y = c | x, W) = p^W(x)_c,$$

$$\begin{aligned} p(Y|X, W) &= \prod_n p^W(x_n)_{y_n=c} \\ &= \prod_n \mathbf{y}_n^T p^W(x_n) \end{aligned}$$

- ▶ prior over W
 - ▶ vectorise W (still write W instead of $\text{vec}(W)$)
 - ▶ same prior as before:

$$p(W) = \mathcal{N}(W; 0_{CK}, s^2 I_{CK})$$

For each c :

$$w_c : k \times 1$$

write

$$W = [w_1, \dots, w_C] \\ : k \times C$$

Def 'logits' =

$$f^W(x) =$$

$$[w_1^T \phi(x), \dots, w_C^T \phi(x)]$$

Def prob func

$$p^W(x) =$$

$$\text{Softmax}(f^W(x))$$

Model:

- ▶ to do predictions

$$p(y^*|x^*, X, Y) = \int p(y^*|x^*, W)p(W|X, Y)dW$$

- ▶ need posterior. But product of softmax and Gaussian is not Gaussian, so can't use tricks from before.. for posterior need evidence:

$$p(Y|X) = \int \prod_n [\mathbf{y}_n^T \text{softmax}(f^W(x_n)_1, \dots, f^W(x_n))] N(W; 0, \mathbf{s}^2 I) dW$$

can't integrate/sum explicitly.. will use VI instead to approx posterior

- ▶ For approx inf need log lik of softmax(f_1, \dots, f_C) = $\left[\frac{e^{f_1}}{e^{f_1} + \dots + e^{f_C}}, \dots\right]$

$$\log p(y = c|x, W) = f_c - \log(e^{f_1} + \dots + e^{f_C})$$

with $[f_1, \dots, f_C]$ the logits vector $[w_1^T \phi(x), \dots, w_C^T \phi(x)]$

- ▶ then expected log likelihood is

$$L(\theta) = \sum_{x_n, y_n=c} \int \left[f^{W(x_n)_c} - \log \left(\sum_{c'} e^{f^{W(x_n)_{c'}}} \right) \right] N(W; \mu_{VI}, \Sigma_{VI}) dW - \text{KL}(q, p)$$

with $f^{W(x)_c} = w_c^T \phi(x)$

- ▶ can't integrate analytically either (log sum exp); need new tools...

- ▶ For approx inf need log lik of softmax(f_1, \dots, f_C) = $\left[\frac{e^{f_1}}{e^{f_1} + \dots + e^{f_C}}, \dots\right]$

$$\log p(y = c|x, W) = f_c - \log(e^{f_1} + \dots + e^{f_C})$$

with $[f_1, \dots, f_C]$ the logits vector $[w_1^T \phi(x), \dots, w_C^T \phi(x)]$

- ▶ then **expected log likelihood** is

$$L(\theta) = \sum_{x_n, y_n=c} \int \left[f^{W(x_n)_c} - \log \left(\sum_{c'} e^{f^{W(x_n)_{c'}}} \right) \right] N(W; \mu_{VI}, \Sigma_{VI}) dW - \text{KL}(q, p)$$

with $f^{W(x)_c} = w_c^T \phi(x)$

- ▶ can't integrate analytically either (log sum exp); need new tools...

- ▶ For approx inf need log lik of softmax(f_1, \dots, f_C) = $\left[\frac{e^{f_1}}{e^{f_1} + \dots + e^{f_C}}, \dots\right]$

$$\log p(y = c|x, W) = f_c - \log(e^{f_1} + \dots + e^{f_C})$$

with $[f_1, \dots, f_C]$ the logits vector $[w_1^T \phi(x), \dots, w_C^T \phi(x)]$

- ▶ then **expected log likelihood** is

$$L(\theta) = \sum_{x_n, y_n=c} \int \left[f^{W(x_n)_c} - \log \left(\sum_{c'} e^{f^{W(x_n)_{c'}}} \right) \right] N(W; \mu_{VI}, \Sigma_{VI}) dW - \text{KL}(q, p)$$

with $f^{W(x)_c} = w_c^T \phi(x)$

- ▶ can't integrate analytically either (log sum exp); need new tools...

Useful tool to estimate expectations

- ▶ let $p(x)$ be some dist which is easy to sample from
- ▶ let $f(x)$ be some function of x
- ▶ assume it to be difficult to eval $E := E_p[f(x)]$
- ▶ can use MC integration instead:
 - ▶ generate $\hat{x}_1, \dots, \hat{x}_T \sim p(x)$
 - ▶ estimate $\hat{E} := 1/T \sum_t f(\hat{x}_t)$
 - ▶ an estimator \hat{E} of E is called *unbiased* if in expectation equals E
 - ▶ \hat{E} is an unbiased estimator of E (prove at home!)

Useful tool to estimate expectations

- ▶ let $p(x)$ be some dist which is easy to sample from
- ▶ let $f(x)$ be some function of x
- ▶ assume it to be difficult to eval $E := E_p[f(x)]$
- ▶ can use MC integration instead:
 - ▶ generate $\hat{x}_1, \dots, \hat{x}_T \sim p(x)$
 - ▶ estimate $\hat{E} := 1/T \sum_t f(\hat{x}_t)$
 - ▶ an estimator \hat{E} of E is called *unbiased* if in expectation equals E
 - ▶ \hat{E} is an unbiased estimator of E (prove at home!)

- ▶ We actually need an estimator of the **derivative of an integral**
- ▶ let $G(\theta)$ be the gradient of $L(\theta)$; will interchangeably use
 - ▶ G (grad of L)
 - ▶ $(L(\theta))' =$ derivative of L wrt θ
 - ▶ $\frac{\partial}{\partial W(\theta)}L(W(\theta)) =$ derivative of L wrt $W(\theta)$
- ▶ if had unbiased derivative estimator $\hat{G}(\theta)$ (estimator of $G(\theta)$) can use a stochastic iterative method to optimise $L(\theta)$:

$$\theta_{n+1} \leftarrow \theta_n + \frac{1}{n} \hat{G}(\theta)$$

go in direction of steepest ascent, on average

- ▶ this is called stochastic gradient descent (well, ascent here)
- ▶ SGD

- ▶ We actually need an estimator of the **derivative of an integral**
- ▶ let $G(\theta)$ be the gradient of $L(\theta)$; will interchangeably use
 - ▶ G (grad of L)
 - ▶ $(L(\theta))' =$ derivative of L wrt θ
 - ▶ $\frac{\partial}{\partial W(\theta)}L(W(\theta)) =$ derivative of L wrt $W(\theta)$
- ▶ if had unbiased derivative estimator $\hat{G}(\theta)$ (estimator of $G(\theta)$) can use a stochastic iterative method to optimise $L(\theta)$:

$$\theta_{n+1} \leftarrow \theta_n + \frac{1}{n} \hat{G}(\theta)$$

go in direction of steepest ascent, on average

- ▶ this is called stochastic gradient descent (well, ascent here)
- ▶ SGD

- ▶ $L(\mu, \sigma) = \int (W + W^2)N(W; \mu, \sigma^2)dW$
 - ▶ can actually eval analytically as $L = \mu + \sigma^2 + \mu^2$
 - ▶ so integral derivative is $G(\mu) = 1 + 2\mu$; will write $G(\mu) := \partial L / \partial \mu$
- ▶ Let's try MC integration first – $\hat{L}(\hat{W}; \mu, \sigma) = \hat{W} + \hat{W}^2$ with realisations (numbers) $\hat{W} \sim \mathcal{N}(\mu, \sigma^2)$, so

$$\hat{G}(\mu) = \partial(\hat{W} + \hat{W}^2) / \partial \mu \stackrel{?}{=} 0$$

(no μ in \hat{L})

- ▶ but \hat{L} clearly depends on μ ;
- ▶ eg increasing μ increases expectation of \hat{L}
- ▶ doesn't look correct... what's going on?

- ▶ $L(\mu, \sigma) = \int (W + W^2) \mathcal{N}(W; \mu, \sigma^2) dW$
 - ▶ can actually eval analytically as $L = \mu + \sigma^2 + \mu^2$
 - ▶ so integral derivative is $G(\mu) = 1 + 2\mu$; will write $G(\mu) := \partial L / \partial \mu$
- ▶ Let's try MC integration first – $\hat{L}(\hat{W}; \mu, \sigma) = \hat{W} + \hat{W}^2$ with realisations (numbers) $\hat{W} \sim \mathcal{N}(\mu, \sigma^2)$, so

$$\hat{G}(\mu) = \partial(\hat{W} + \hat{W}^2) / \partial \mu \stackrel{?}{=} 0$$

(no μ in \hat{L})

- ▶ but \hat{L} clearly depends on μ ;
- ▶ eg increasing μ increases expectation of \hat{L}
- ▶ doesn't look correct... what's going on?

- ▶ \hat{L} deps on μ through \hat{W} ; \hat{W} is actually a function of μ as well as a rv $\hat{\epsilon}$ indep of θ :

$$\hat{W} \sim \mathcal{N}(\mu, \sigma^2) \quad \leftrightarrow \quad \hat{W} = \hat{W}(\theta, \hat{\epsilon}) = \mu + \sigma\hat{\epsilon}; \quad \hat{\epsilon} \sim \mathcal{N}(0, 1)$$

- ▶ then can rewrite \hat{L} as

$$\hat{L}(\mu, \sigma, \hat{\epsilon}) = (\mu + \sigma\hat{\epsilon}) + (\mu + \sigma\hat{\epsilon})^2 = \mu + \mu^2 + 2\mu\sigma\hat{\epsilon} + \sigma\hat{\epsilon} + \sigma^2\hat{\epsilon}^2$$

and

$$\hat{G}(\mu) = 1 + 2\mu + 2\sigma\hat{\epsilon}$$

- ▶ check:

$$E_{\rho(\epsilon)}[\hat{G}] = 1 + 2\mu = G$$

ie \hat{G} is an unbiased estimator of G

- ▶ \hat{L} deps on μ through \hat{W} ; \hat{W} is actually a function of μ as well as a rv $\hat{\epsilon}$ indep of θ :

$$\hat{W} \sim \mathcal{N}(\mu, \sigma^2) \quad \leftrightarrow \quad \hat{W} = \hat{W}(\theta, \hat{\epsilon}) = \mu + \sigma\hat{\epsilon}; \quad \hat{\epsilon} \sim \mathcal{N}(0, 1)$$

- ▶ then can rewrite \hat{L} as

$$\hat{L}(\mu, \sigma, \hat{\epsilon}) = (\mu + \sigma\hat{\epsilon}) + (\mu + \sigma\hat{\epsilon})^2 = \mu + \mu^2 + 2\mu\sigma\hat{\epsilon} + \sigma\hat{\epsilon} + \sigma^2\hat{\epsilon}^2$$

and

$$\hat{G}(\mu) = 1 + 2\mu + 2\sigma\hat{\epsilon}$$

- ▶ check:

$$E_{\rho(\epsilon)}[\hat{G}] = 1 + 2\mu = G$$

ie \hat{G} is an unbiased estimator of G

- ▶ \hat{L} deps on μ through \hat{W} ; \hat{W} is actually a function of μ as well as a rv $\hat{\epsilon}$ indep of θ :

$$\hat{W} \sim \mathcal{N}(\mu, \sigma^2) \quad \Leftrightarrow \quad \hat{W} = \hat{W}(\theta, \hat{\epsilon}) = \mu + \sigma\hat{\epsilon}; \quad \hat{\epsilon} \sim \mathcal{N}(0, 1)$$

- ▶ then can rewrite \hat{L} as

$$\hat{L}(\mu, \sigma, \hat{\epsilon}) = (\mu + \sigma\hat{\epsilon}) + (\mu + \sigma\hat{\epsilon})^2 = \mu + \mu^2 + 2\mu\sigma\hat{\epsilon} + \sigma\hat{\epsilon} + \sigma^2\hat{\epsilon}^2$$

and

$$\hat{G}(\mu) = 1 + 2\mu + 2\sigma\hat{\epsilon}$$

- ▶ check:

$$E_{\rho(\epsilon)}[\hat{G}] = 1 + 2\mu = G$$

ie \hat{G} is an unbiased estimator of G

- ▶ technique known in literature as the *re-parametrisation trick*
 - ▶ also known as a *pathwise derivative estimator*, *infinitesimal perturbation analysis*, and *stochastic backpropagation*
- ▶ in general:
 - ▶ given func $f(W)$, dist $q_\theta(W)$
 - ▶ want to estimate gradients of $L(\theta) = \int f(W)q_\theta(W)dW$
 - ▶ if W can be reparam as $W = g(\theta, \epsilon)$ with ϵ not dependent on θ , and g is differentiable wrt θ
 - ▶ then $\hat{G}(\theta, \epsilon) = f'(g(\theta, \epsilon)) \cdot \partial g(\theta, \epsilon) / \partial \theta$
 - ▶ .. and plug into a stochastic optimiser
- ▶ eg, for Gaussian q ...
 - ▶ $W = g([\mu, \sigma], \epsilon) = \mu + \sigma\epsilon$
 - ▶ $\partial g([\mu, \sigma], \epsilon) / \partial \mu = 1$ and $\partial g([\mu, \sigma], \epsilon) / \partial \sigma = \epsilon$
 - ▶ so $\hat{G}(\epsilon; \mu) = f'(\mu + \sigma\epsilon) \cdot 1$ and $\hat{G}(\epsilon; \sigma) = f'(\mu + \sigma\epsilon)\epsilon$
 - ▶ with $\epsilon \sim \mathcal{N}(0, 1)$
 - ▶ can substitute in $\hat{W} = \mu + \sigma\epsilon$: sample $\hat{W} \sim q_\theta(W)$; then $\hat{G}(\hat{W}; \mu) = f'(\hat{W})$ and $\hat{G}(\hat{W}; \sigma) = f'(\hat{W})(\hat{W} - \mu) / \sigma$.

- ▶ technique known in literature as the *re-parametrisation trick*
 - ▶ also known as a *pathwise derivative estimator*, *infinitesimal perturbation analysis*, and *stochastic backpropagation*
- ▶ in general:
 - ▶ given func $f(W)$, dist $q_\theta(W)$
 - ▶ want to estimate gradients of $L(\theta) = \int f(W)q_\theta(W)dW$
 - ▶ if W can be reparam as $W = g(\theta, \epsilon)$ with ϵ not dependent on θ , and g is **differentiable wrt θ**
 - ▶ then $\hat{G}(\theta, \hat{\epsilon}) = f'(g(\theta, \hat{\epsilon})) \cdot \partial g(\theta, \hat{\epsilon})/\partial \theta$
 - ▶ .. and plug into a stochastic optimiser
- ▶ eg, for Gaussian $q...$
 - ▶ $W = g([\mu, \sigma], \epsilon) = \mu + \sigma\epsilon$
 - ▶ $\partial g([\mu, \sigma], \epsilon)/\partial \mu = 1$ and $\partial g([\mu, \sigma], \epsilon)/\partial \sigma = \epsilon$
 - ▶ so $\hat{G}(\hat{\epsilon}; \mu) = f'(\mu + \sigma\hat{\epsilon}) \cdot 1$ and $\hat{G}(\hat{\epsilon}; \sigma) = f'(\mu + \sigma\hat{\epsilon})\hat{\epsilon}$
 - ▶ with $\hat{\epsilon} \sim \mathcal{N}(0, 1)$
 - ▶ can substitute in $\hat{W} = \mu + \sigma\hat{\epsilon}$: sample $\hat{W} \sim q_\theta(W)$; then $\hat{G}(\hat{W}; \mu) = f'(\hat{W})$ and $\hat{G}(\hat{W}; \sigma) = f'(\hat{W})(\hat{W} - \mu)/\sigma$.

- ▶ technique known in literature as the *re-parametrisation trick*
 - ▶ also known as a *pathwise derivative estimator*, *infinitesimal perturbation analysis*, and *stochastic backpropagation*
- ▶ in general:
 - ▶ given func $f(W)$, dist $q_\theta(W)$
 - ▶ want to estimate gradients of $L(\theta) = \int f(W)q_\theta(W)dW$
 - ▶ if W can be reparam as $W = g(\theta, \epsilon)$ with ϵ not dependent on θ , and g is **differentiable wrt θ**
 - ▶ then $\hat{G}(\theta, \hat{\epsilon}) = f'(g(\theta, \hat{\epsilon})) \cdot \partial g(\theta, \hat{\epsilon})/\partial \theta$
 - ▶ .. and plug into a stochastic optimiser
- ▶ eg, for Gaussian $q...$
 - ▶ $W = g([\mu, \sigma], \epsilon) = \mu + \sigma\epsilon$
 - ▶ $\partial g([\mu, \sigma], \epsilon)/\partial \mu = 1$ and $\partial g([\mu, \sigma], \epsilon)/\partial \sigma = \epsilon$
 - ▶ so $\hat{G}(\hat{\epsilon}; \mu) = f'(\mu + \sigma\hat{\epsilon}) \cdot 1$ and $\hat{G}(\hat{\epsilon}; \sigma) = f'(\mu + \sigma\hat{\epsilon})\hat{\epsilon}$
 - ▶ with $\hat{\epsilon} \sim \mathcal{N}(0, 1)$
 - ▶ can substitute in $\hat{W} = \mu + \sigma\hat{\epsilon}$: sample $\hat{W} \sim q_\theta(W)$; then $\hat{G}(\hat{W}; \mu) = f'(\hat{W})$ and $\hat{G}(\hat{W}; \sigma) = f'(\hat{W})(\hat{W} - \mu)/\sigma$.

- ▶ technique known in literature as the *re-parametrisation trick*
 - ▶ also known as a *pathwise derivative estimator*, *infinitesimal perturbation analysis*, and *stochastic backpropagation*
- ▶ in general:
 - ▶ given func $f(W)$, dist $q_\theta(W)$
 - ▶ want to estimate gradients of $L(\theta) = \int f(W)q_\theta(W)dW$
 - ▶ if W can be reparam as $W = g(\theta, \epsilon)$ with ϵ not dependent on θ , and g is **differentiable wrt θ**
 - ▶ then $\hat{G}(\theta, \hat{\epsilon}) = f'(g(\theta, \hat{\epsilon})) \cdot \partial g(\theta, \hat{\epsilon}) / \partial \theta$
 - ▶ .. and plug into a stochastic optimiser
- ▶ eg, for Gaussian $q...$
 - ▶ $W = g([\mu, \sigma], \epsilon) = \mu + \sigma\epsilon$
 - ▶ $\partial g([\mu, \sigma], \epsilon) / \partial \mu = 1$ and $\partial g([\mu, \sigma], \epsilon) / \partial \sigma = \epsilon$
 - ▶ so $\hat{G}(\epsilon; \mu) = f'(\mu + \sigma\epsilon) \cdot 1$ and $\hat{G}(\epsilon; \sigma) = f'(\mu + \sigma\epsilon)\epsilon$
 - ▶ with $\hat{\epsilon} \sim \mathcal{N}(0, 1)$
 - ▶ can substitute in $\hat{W} = \mu + \sigma\hat{\epsilon}$: sample $\hat{W} \sim q_\theta(W)$; then $\hat{G}(W; \mu) = f'(\hat{W})$ and $\hat{G}(W; \sigma) = f'(\hat{W})(\hat{W} - \mu) / \sigma$.

- ▶ technique known in literature as the *re-parametrisation trick*
 - ▶ also known as a *pathwise derivative estimator*, *infinitesimal perturbation analysis*, and *stochastic backpropagation*
- ▶ in general:
 - ▶ given func $f(W)$, dist $q_\theta(W)$
 - ▶ want to estimate gradients of $L(\theta) = \int f(W)q_\theta(W)dW$
 - ▶ if W can be reparam as $W = g(\theta, \epsilon)$ with ϵ not dependent on θ , and g is **differentiable wrt θ**
 - ▶ then $\hat{G}(\theta, \hat{\epsilon}) = f'(g(\theta, \hat{\epsilon})) \cdot \partial g(\theta, \hat{\epsilon})/\partial \theta$
 - ▶ .. and plug into a stochastic optimiser
- ▶ eg, for Gaussian $q...$
 - ▶ $W = g([\mu, \sigma], \epsilon) = \mu + \sigma \epsilon$
 - ▶ $\partial g([\mu, \sigma], \epsilon)/\partial \mu = \mathbf{1}$ and $\partial g([\mu, \sigma], \epsilon)/\partial \sigma = \epsilon$
 - ▶ so $\hat{G}(\hat{\epsilon}; \mu) = f'(\mu + \sigma \hat{\epsilon}) \cdot \mathbf{1}$ and $\hat{G}(\hat{\epsilon}; \sigma) = f'(\mu + \sigma \hat{\epsilon}) \hat{\epsilon}$
 - ▶ with $\hat{\epsilon} \sim \mathcal{N}(0, I)$
 - ▶ can substitute in $\hat{W} = \mu + \sigma \hat{\epsilon}$: sample $\hat{W} \sim q_\theta(W)$; then $\hat{G}(\hat{W}; \mu) = f'(\hat{W})$ and $\hat{G}(\hat{W}; \sigma) = f'(\hat{W})(\hat{W} - \mu)/\sigma$.

- ▶ technique known in literature as the *re-parametrisation trick*
 - ▶ also known as a *pathwise derivative estimator*, *infinitesimal perturbation analysis*, and *stochastic backpropagation*
- ▶ in general:
 - ▶ given func $f(W)$, dist $q_\theta(W)$
 - ▶ want to estimate gradients of $L(\theta) = \int f(W)q_\theta(W)dW$
 - ▶ if W can be reparam as $W = g(\theta, \epsilon)$ with ϵ not dependent on θ , and g is **differentiable wrt** θ
 - ▶ then $\hat{G}(\theta, \hat{\epsilon}) = f'(g(\theta, \hat{\epsilon})) \cdot \partial g(\theta, \hat{\epsilon})/\partial \theta$
 - ▶ .. and plug into a stochastic optimiser
- ▶ eg, for Gaussian q ...
 - ▶ $W = g([\mu, \sigma], \epsilon) = \mu + \sigma \epsilon$
 - ▶ $\partial g([\mu, \sigma], \epsilon)/\partial \mu = \mathbf{1}$ and $\partial g([\mu, \sigma], \epsilon)/\partial \sigma = \epsilon$
 - ▶ so $\hat{G}(\hat{\epsilon}; \mu) = f'(\mu + \sigma \hat{\epsilon}) \cdot \mathbf{1}$ and $\hat{G}(\hat{\epsilon}; \sigma) = f'(\mu + \sigma \hat{\epsilon}) \hat{\epsilon}$
 - ▶ with $\hat{\epsilon} \sim \mathcal{N}(0, I)$
 - ▶ can substitute in $\hat{W} = \mu + \sigma \hat{\epsilon}$: sample $\hat{W} \sim q_\theta(W)$; then $\hat{G}(\hat{W}; \mu) = f'(\hat{W})$ and $\hat{G}(\hat{W}; \sigma) = f'(\hat{W})(\hat{W} - \mu)/\sigma$.

- ▶ technique known in literature as the *re-parametrisation trick*
 - ▶ also known as a *pathwise derivative estimator*, *infinitesimal perturbation analysis*, and *stochastic backpropagation*
- ▶ in general:
 - ▶ given func $f(W)$, dist $q_\theta(W)$
 - ▶ want to estimate gradients of $L(\theta) = \int f(W)q_\theta(W)dW$
 - ▶ if W can be reparam as $W = g(\theta, \epsilon)$ with ϵ not dependent on θ , and g is **differentiable wrt θ**
 - ▶ then $\hat{G}(\theta, \hat{\epsilon}) = f'(g(\theta, \hat{\epsilon})) \cdot \partial g(\theta, \hat{\epsilon})/\partial \theta$
 - ▶ .. and plug into a stochastic optimiser
- ▶ eg, for Gaussian $q...$
 - ▶ $W = g([\mu, \sigma], \epsilon) = \mu + \sigma \epsilon$
 - ▶ $\partial g([\mu, \sigma], \epsilon)/\partial \mu = 1$ and $\partial g([\mu, \sigma], \epsilon)/\partial \sigma = \epsilon$
 - ▶ so $\hat{G}(\hat{\epsilon}; \mu) = f'(\mu + \sigma \hat{\epsilon}) \cdot 1$ and $\hat{G}(\hat{\epsilon}; \sigma) = f'(\mu + \sigma \hat{\epsilon}) \hat{\epsilon}$
 - ▶ with $\hat{\epsilon} \sim \mathcal{N}(0, I)$
 - ▶ can substitute in $\hat{W} = \mu + \sigma \hat{\epsilon}$: sample $\hat{W} \sim q_\theta(W)$; then $\hat{G}(\hat{W}; \mu) = f'(\hat{W})$ and $\hat{G}(\hat{W}; \sigma) = f'(\hat{W})(\hat{W} - \mu)/\sigma$.

Remember prev ELBO which we couldn't eval

$$L(\theta) = \sum_{x_n, y_n=c} \int \left[f^W(x_n)_c - \log \left(\sum_{c'} e^{f^W(x_n)_{c'}} \right) \right] N(W; \mu_{V1}, \Sigma_{V1}) dW - \text{KL}(q, p)$$

W vectorised w dim CK by 1, so is μ_{V1} , and assume Σ_{V1} is diagonal w dim CK by CK

- ▶ using MC integration
 - ▶ sample $\hat{\epsilon} \sim \mathcal{N}(0, I_{CK})$
 - ▶ write $\text{vec} \hat{W}(\theta, \hat{\epsilon}) = \mu_{V1} + \Sigma_{V1}^{1/2} \hat{\epsilon}$
 - ▶ reshape $\text{vec} \hat{W}$ to K by C : $\hat{W}(\theta, \hat{\epsilon})$
 - ▶ write $f^{\theta, \hat{\epsilon}}(x) = f^{\hat{W}(\theta, \hat{\epsilon})}(x)$
 - ▶ giving

$$\hat{L}(\theta, \hat{\epsilon}) = \sum_{x_n, y_n=c} f^{\theta, \hat{\epsilon}}(x_n)_c - \log \left(\sum_{c'} e^{f^{\theta, \hat{\epsilon}}(x_n)_{c'}} \right) - \text{KL}(q, p)$$

- ▶ with $E_{p(\epsilon)}[\hat{L}(\theta, \epsilon)] = L(\theta)$, $E_{p(\epsilon)}[\hat{G}(\theta, \epsilon)] = G(\theta)$

Remember prev ELBO which we couldn't eval

$$L(\theta) = \sum_{x_n, y_n=c} \int \left[f^W(x_n)_c - \log \left(\sum_{c'} e^{f^W(x_n)_{c'}} \right) \right] N(W; \mu_{V1}, \Sigma_{V1}) dW - \text{KL}(q, p)$$

W vectorised w dim CK by 1, so is μ_{V1} , and assume Σ_{V1} is diagonal w dim CK by CK

- ▶ using MC integration
 - ▶ sample $\hat{\epsilon} \sim \mathcal{N}(0, I_{CK})$
 - ▶ write $\text{vec} \hat{W}(\theta, \hat{\epsilon}) = \mu_{V1} + \Sigma_{V1}^{1/2} \hat{\epsilon}$
 - ▶ reshape $\text{vec} \hat{W}$ to K by C : $\hat{W}(\theta, \hat{\epsilon})$
 - ▶ write $f^{\theta, \hat{\epsilon}}(x) = f^{\hat{W}(\theta, \hat{\epsilon})}(x)$
 - ▶ giving

$$\hat{L}(\theta, \hat{\epsilon}) = \sum_{x_n, y_n=c} f^{\theta, \hat{\epsilon}}(x_n)_c - \log \left(\sum_{c'} e^{f^{\theta, \hat{\epsilon}}(x_n)_{c'}} \right) - \text{KL}(q, p)$$

- ▶ with $E_{p(\epsilon)}[\hat{L}(\theta, \epsilon)] = L(\theta)$, $E_{p(\epsilon)}[\hat{G}(\theta, \epsilon)] = G(\theta)$

Remember prev ELBO which we couldn't eval

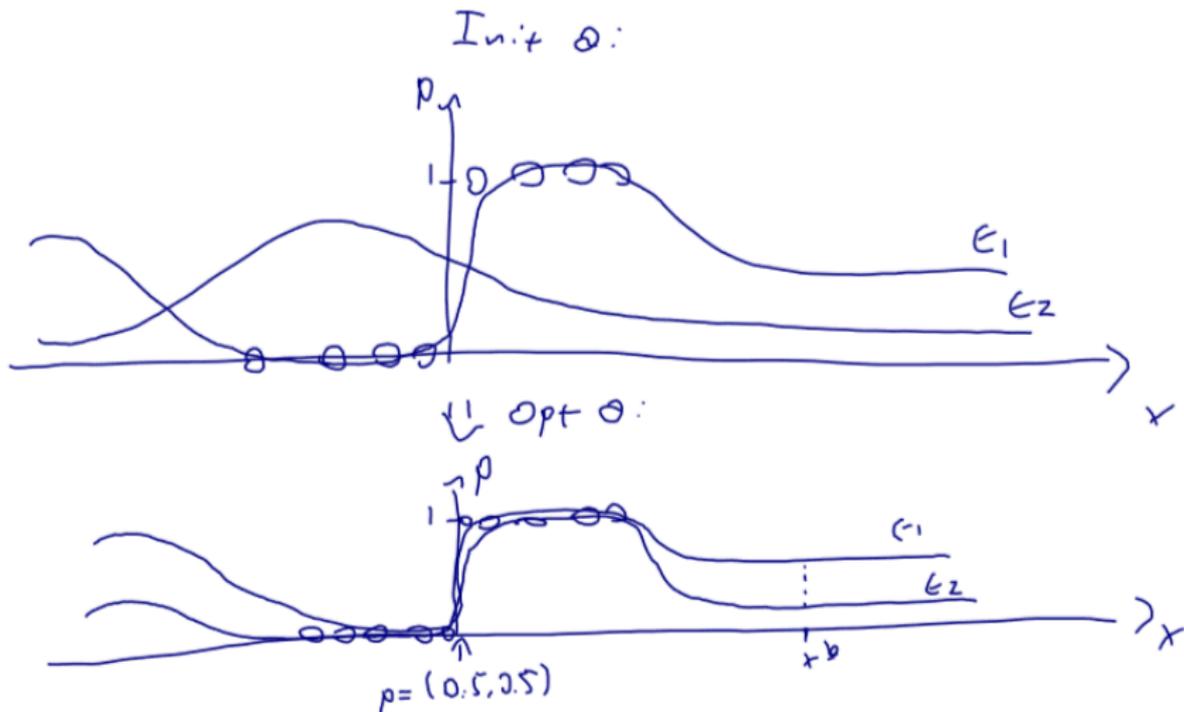
$$L(\theta) = \sum_{x_n, y_n=c} \int \left[f^W(x_n)_c - \log \left(\sum_{c'} e^{f^W(x_n)_{c'}} \right) \right] N(W; \mu_{V1}, \Sigma_{V1}) dW - \text{KL}(q, p)$$

W vectorised w dim CK by 1, so is μ_{V1} , and assume Σ_{V1} is diagonal w dim CK by CK

- ▶ using MC integration
 - ▶ sample $\hat{\epsilon} \sim \mathcal{N}(0, I_{CK})$
 - ▶ write $\text{vec} \hat{W}(\theta, \hat{\epsilon}) = \mu_{V1} + \Sigma_{V1}^{1/2} \hat{\epsilon}$
 - ▶ reshape $\text{vec} \hat{W}$ to K by C : $\hat{W}(\theta, \hat{\epsilon})$
 - ▶ write $f^{\theta, \hat{\epsilon}}(x) = f^{\hat{W}(\theta, \hat{\epsilon})}(x)$
 - ▶ giving

$$\hat{L}(\theta, \hat{\epsilon}) = \sum_{x_n, y_n=c} f^{\theta, \hat{\epsilon}}(x_n)_c - \log \left(\sum_{c'} e^{f^{\theta, \hat{\epsilon}}(x_n)_{c'}} \right) - \text{KL}(q, p)$$

- ▶ with $E_{p(\epsilon)}[\hat{L}(\theta, \epsilon)] = L(\theta)$, $E_{p(\epsilon)}[\hat{G}(\theta, \epsilon)] = G(\theta)$



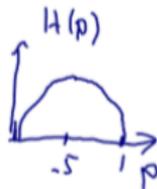
Uncertainty in Classification

Epistemic uncertainty in classification (vs regression)

- ▶ finally have tools to get epistemic uncertainty for classification
- ▶ but quantifying uncertainty in classification is not as straightforward as in regression...
- ▶ use various *measures of uncertainty* from the field of Information Theory, which have different properties
- ▶ each capturing different uncertainty desiderata

Useful tools

- ▶ Entropy $H_{p(X)}[X] = -\sum_{\text{outcomes } x} p(X=x) \log p(X=x)$
 - ▶ high when p is **uniform**, 0 when one outcome is **certain**
- ▶ Mutual information of rvs X and Y



$$MI(X, Y) = H_{p(X)}[X] - E_{p(Y)}[H_{p(X|Y)}[X]]$$

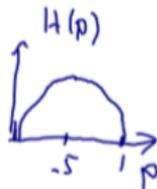
- ▶ “how much information on X we would get if we had observed Y ”

Epistemic uncertainty in classification (vs regression)

- ▶ finally have tools to get epistemic uncertainty for classification
- ▶ but quantifying uncertainty in classification is not as straightforward as in regression...
- ▶ use various *measures of uncertainty* from the field of Information Theory, which have different properties
- ▶ each capturing different uncertainty desiderata

Useful tools

- ▶ Entropy $H_{p(x)}[X] = -\sum_{\text{outcomes } x} p(X=x) \log p(X=x)$
 - ▶ high when p is **uniform**, 0 when one outcome is **certain**
- ▶ Mutual information of rvs X and Y



$$MI(X, Y) = H_{p(x)}[X] - E_{p(y)}[H_{p(x|y)}[X]]$$

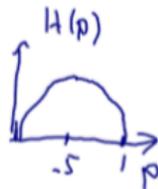
- ▶ “how much information on X we would get if we had observed Y ”

Epistemic uncertainty in classification (vs regression)

- ▶ finally have tools to get epistemic uncertainty for classification
- ▶ but quantifying uncertainty in classification is not as straightforward as in regression...
- ▶ use various *measures of uncertainty* from the field of Information Theory, which have different properties
- ▶ each capturing different uncertainty desiderata

Useful tools

- ▶ Entropy $H_{p(X)}[X] = -\sum_{\text{outcomes } x} p(X=x) \log p(X=x)$
 - ▶ high when p is **uniform**, 0 when one outcome is **certain**
- ▶ Mutual information of rvs X and Y



$$MI(X, Y) = H_{p(X)}[X] - E_{p(Y)}[H_{p(X|Y)}[X]]$$

- ▶ “how much information on X we would get if we had observed Y ”

A quick overview:

▶ *Predictive Entropy*

- ▶ entropy of predictive distribution $p(y = y^* | x^*, \mathcal{D})$

$$H_{p(y^* | x^*, \mathcal{D})}[y^*] = - \sum_{y^*=c} p(y^* = c | x^*, \mathcal{D}) \log p(y^* = c | x^*, \mathcal{D})$$

▶ *Mutual Information (MI)*

- ▶ between model params rv W and model output rv y^* on input x^*

$$MI(y^*, W | \mathcal{D}, x^*) = H_{p(y^* | x^*, \mathcal{D})}[y^*] - E_{p(W | \mathcal{D})}[H_{p(y^* | x^*, W)}[y^*]]$$

- ▶ satisfies

$$0 \leq MI[x^*] \leq H[x^*]$$

A quick overview:

▶ *Predictive Entropy*

- ▶ entropy of predictive distribution $p(y = y^* | x^*, \mathcal{D})$

$$H_{p(y^* | x^*, \mathcal{D})}[y^*] = - \sum_{y^*=c} p(y^* = c | x^*, \mathcal{D}) \log p(y^* = c | x^*, \mathcal{D})$$

▶ *Mutual Information (MI)*

- ▶ between model params rv W and model output rv y^* on input x^*

$$MI(y^*, W | \mathcal{D}, x^*) = H_{p(y^* | x^*, \mathcal{D})}[y^*] - E_{p(W | \mathcal{D})}[H_{p(y^* | x^*, W)}[y^*]]$$

- ▶ satisfies

$$0 \leq MI[x^*] \leq H[x^*]$$

Predictive entropy

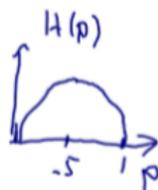
$$H_{p(y^*|x^*, \mathcal{D})}[y^*] = - \sum_{y^*=c} p(y^* = c|x^*, \mathcal{D}) \log p(y^* = c|x^*, \mathcal{D})$$

- ▶ MC approximation

$$p(y^* = c|x^*, \mathcal{D}) \approx \frac{1}{T} \sum_t p^{\hat{W}_t}(x^*)_c$$

with $\hat{W}_t \sim q_\theta(W)$ and $p^{\hat{W}_t}(x^*) = \text{softmax}(f^{\hat{W}_t}(x^*))$

- ▶ high when predictive is near uniform
- ▶ so, high **either** when we have inherent ambiguity
 - ▶ eg when a point x has training labels both 0 and 1
 - ▶ for ambiguous input x loss is $\log p(x) + \log(1 - p(x))$
 - ▶ cross entropy loss minimiser (=ELBO maximiser) is to predict $p = .5$
 - ▶ all func draws will go through $(.5, .5)$ (ie high entropy)
- ▶ or when far away from data: eg half draws=1 and half draws=0



Predictive entropy

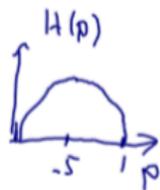
$$H_{p(y^*|x^*, \mathcal{D})}[y^*] = - \sum_{y^*=c} p(y^* = c|x^*, \mathcal{D}) \log p(y^* = c|x^*, \mathcal{D})$$

- ▶ MC approximation

$$p(y^* = c|x^*, \mathcal{D}) \approx \frac{1}{T} \sum_t p^{\hat{W}_t}(x^*)_c$$

with $\hat{W}_t \sim q_\theta(W)$ and $p^{\hat{W}_t}(x^*) = \text{softmax}(f^{\hat{W}_t}(x^*))$

- ▶ high when predictive is near uniform
- ▶ so, high **either** when we have inherent ambiguity
 - ▶ eg when a point x has training labels both 0 and 1
 - ▶ for ambiguous input x loss is $\log p(x) + \log(1 - p(x))$
 - ▶ cross entropy loss minimiser (=ELBO maximiser) is to predict $p = .5$
 - ▶ all func draws will go through $(.5, .5)$ (ie high entropy)
- ▶ or when far away from data: eg half draws=1 and half draws=0



Predictive entropy

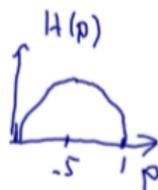
$$H_{p(y^*|x^*, \mathcal{D})}[y^*] = - \sum_{y^*=c} p(y^* = c|x^*, \mathcal{D}) \log p(y^* = c|x^*, \mathcal{D})$$

- ▶ MC approximation

$$p(y^* = c|x^*, \mathcal{D}) \approx \frac{1}{T} \sum_t p^{\hat{W}_t}(x^*)_c$$

with $\hat{W}_t \sim q_\theta(W)$ and $p^{\hat{W}_t}(x^*) = \text{softmax}(f^{\hat{W}_t}(x^*))$

- ▶ high when predictive is near uniform
- ▶ so, high **either** when we have inherent ambiguity
 - ▶ eg when a point x has training labels both 0 and 1
 - ▶ for ambiguous input x loss is $\log p(x) + \log(1 - p(x))$
 - ▶ cross entropy loss minimiser (=ELBO maximiser) is to predict $p = .5$
 - ▶ all func draws will go through $(.5, .5)$ (ie high entropy)
- ▶ **or** when far away from data: eg half draws=1 and half draws=0



Mutual information

$$MI(y^*, W|\mathcal{D}, x^*) = H_{p(y^*|x^*, \mathcal{D})}[y^*] - E_{p(W|\mathcal{D})}[H_{p(y^*|x^*, W)}[y^*]]$$

- ▶ MI MC approx (second term)

$$\begin{aligned} & \int p(W|\mathcal{D}) \sum_{y^*=c} p(y^* = c|x^*, W) \log p(y^* = c|x^*, W) dW \\ & \approx \frac{1}{T} \sum_{t, y^*=c} p^{\hat{W}_t}(x^*)_c \log p^{\hat{W}_t}(x^*)_c \end{aligned}$$

with $\hat{W}_t \sim q_\theta(W)$

- ▶ high **only** when we are far away from data
 - ▶ has “second term = first term” if all func draws same for input x
 - ▶ “second term = 0” when func preds are confident and all over the place
- ▶ ie, capturing **only** epistemic uncertainty (vs pred ent capturing epistemic *and* aleatoric uncertainty)

Mutual information

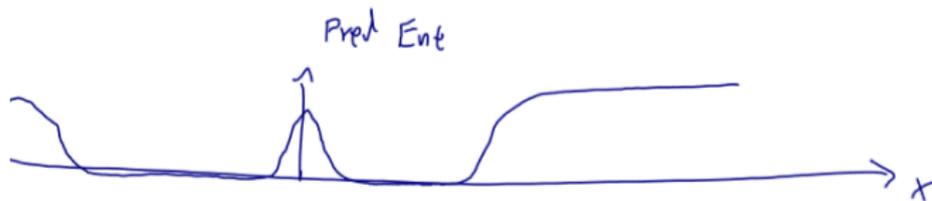
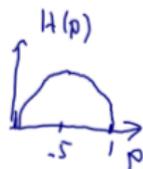
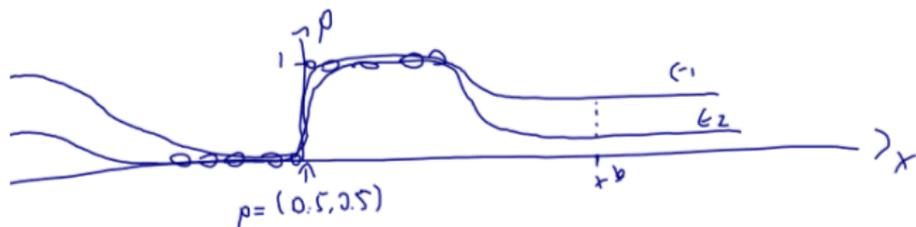
$$MI(y^*, W|\mathcal{D}, x^*) = H_{p(y^*|x^*, \mathcal{D})}[y^*] - E_{p(W|\mathcal{D})}[H_{p(y^*|x^*, W)}[y^*]]$$

- ▶ MI MC approx (second term)

$$\begin{aligned} & \int p(W|\mathcal{D}) \sum_{y^*=c} p(y^* = c|x^*, W) \log p(y^* = c|x^*, W) dW \\ & \approx \frac{1}{T} \sum_{t, y^*=c} p^{\hat{W}_t}(x^*)_c \log p^{\hat{W}_t}(x^*)_c \end{aligned}$$

with $\hat{W}_t \sim q_\theta(W)$

- ▶ high **only** when we are far away from data
 - ▶ has “second term = first term” if all func draws same for input x
 - ▶ “second term = 0” when func preds are confident and all over the place
- ▶ ie, capturing **only** epistemic uncertainty (vs pred ent capturing epistemic *and* aleatoric uncertainty)



Predictive: $p(y^* = c | x^*, \mathcal{D}) \approx \frac{1}{T} \sum_t p^{\hat{W}_t}(x^*)_c$

MI: $MI(y^*, W | \mathcal{D}, x^*) = H_{p(y^* | x^*, \mathcal{D})}[y^*] - E_{p(W | \mathcal{D})}[H_{p(y^* | x^*, W)}[y^*]]$

Stochastic Approximate Inference in Deep NN

- ▶ Model for regression (D outputs) / classification
- ▶ ELBO $L(\theta) = \int q_\theta(W) \log p(Y|X, W) dW - \text{KL}(q, \text{prior})$
- ▶ log likelihood eg

$$\log p(Y|X, W) = -\frac{1}{2\sigma^2} \sum \|y_n - f^W(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2$$
- ▶ approx post eg $q_\theta(w_{kd}) = N(w_{kd}; m_{kd}, \sigma_{kd}^2)$
 - ▶ $\text{KL}(q, \text{prior}) = \sum_{kd} 1/2(s^{-2}\sigma_{kd}^2 + s^{-2}m_{kd}^2 - 1 + \log(s^2/\sigma_{kd}^2))$
- ▶ MC integration:
 - ▶ sample $\hat{\epsilon}_{kd} \sim \mathcal{N}(0, 1)$ and write $\hat{w}_{kd} = m_{kd} + \sigma_{kd}\hat{\epsilon}_{kd}$, $\hat{\epsilon} = \{\hat{\epsilon}_{kd}\}$
 - ▶ giving a K by D stochastic weight matrix: $\hat{W}(\theta, \hat{\epsilon})$
 - ▶ write $f^{\theta, \hat{\epsilon}}(x) = \hat{W}(\theta, \hat{\epsilon})^T \phi(x)$
 - ▶ giving

$$\hat{L}(\theta, \{\hat{\epsilon}_n\}) = -\frac{1}{2\sigma^2} \sum_{x_n, y_n} \|y_n - f^{\theta, \hat{\epsilon}_n}(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2s^2} \|M\|_2^2 \dots$$

- ▶ Model for regression (D outputs) / classification
- ▶ ELBO $L(\theta) = \int q_\theta(W) \log p(Y|X, W) dW - \text{KL}(q, \text{prior})$
- ▶ log likelihood eg

$$\log p(Y|X, W) = -\frac{1}{2\sigma^2} \sum \|y_n - f^W(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2$$
- ▶ approx post eg $q_\theta(w_{kd}) = N(w_{kd}; m_{kd}, \sigma_{kd}^2)$
 - ▶ $\text{KL}(q, \text{prior}) = \sum_{kd} 1/2(s^{-2}\sigma_{kd}^2 + s^{-2}m_{kd}^2 - 1 + \log(s^2/\sigma_{kd}^2))$
- ▶ MC integration:
 - ▶ sample $\hat{\epsilon}_{kd} \sim \mathcal{N}(0, 1)$ and write $\hat{w}_{kd} = m_{kd} + \sigma_{kd}\hat{\epsilon}_{kd}$, $\hat{\epsilon} = \{\hat{\epsilon}_{kd}\}$
 - ▶ giving a K by D stochastic weight matrix: $\hat{W}(\theta, \hat{\epsilon})$
 - ▶ write $f^{\theta, \hat{\epsilon}}(x) = \hat{W}(\theta, \hat{\epsilon})^T \phi(x)$
 - ▶ giving

$$\hat{L}(\theta, \{\hat{\epsilon}_n\}) = -\frac{1}{2\sigma^2} \sum_{x_n, y_n} \|y_n - f^{\theta, \hat{\epsilon}_n}(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2s^2} \|M\|_2^2 \dots$$

- ▶ Model for regression (D outputs) / classification
- ▶ ELBO $L(\theta) = \int q_\theta(W) \log p(Y|X, W) dW - \text{KL}(q, \text{prior})$
- ▶ log likelihood eg

$$\log p(Y|X, W) = -\frac{1}{2\sigma^2} \sum \|y_n - f^W(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2$$
- ▶ approx post eg $q_\theta(w_{kd}) = N(w_{kd}; m_{kd}, \sigma_{kd}^2)$
 - ▶ $\text{KL}(q, \text{prior}) = \sum_{kd} 1/2(s^{-2}\sigma_{kd}^2 + s^{-2}m_{kd}^2 - 1 + \log(s^2/\sigma_{kd}^2))$
- ▶ MC integration:
 - ▶ sample $\hat{\epsilon}_{kd} \sim \mathcal{N}(0, 1)$ and write $\hat{w}_{kd} = m_{kd} + \sigma_{kd}\hat{\epsilon}_{kd}$, $\hat{\epsilon} = \{\hat{\epsilon}_{kd}\}$
 - ▶ giving a K by D stochastic weight matrix: $\hat{W}(\theta, \hat{\epsilon})$
 - ▶ write $f^{\theta, \hat{\epsilon}}(x) = \hat{W}(\theta, \hat{\epsilon})^T \phi(x)$
 - ▶ giving

$$\hat{L}(\theta, \{\hat{\epsilon}_n\}) = -\frac{1}{2\sigma^2} \sum_{x_n, y_n} \|y_n - f^{\theta, \hat{\epsilon}_n}(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2s^2} \|M\|_2^2 \dots$$

- ▶ Model for regression (D outputs) / classification
- ▶ ELBO $L(\theta) = \int q_\theta(W) \log p(Y|X, W) dW - \text{KL}(q, \text{prior})$
- ▶ log likelihood eg

$$\log p(Y|X, W) = -\frac{1}{2\sigma^2} \sum \|y_n - f^W(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2$$
- ▶ approx post eg $q_\theta(w_{kd}) = N(w_{kd}; m_{kd}, \sigma_{kd}^2)$
 - ▶ $\text{KL}(q, \text{prior}) = \sum_{kd} 1/2(s^{-2}\sigma_{kd}^2 + s^{-2}m_{kd}^2 - 1 + \log(s^2/\sigma_{kd}^2))$
- ▶ MC integration:
 - ▶ sample $\hat{\epsilon}_{kd} \sim \mathcal{N}(0, 1)$ and write $\hat{w}_{kd} = m_{kd} + \sigma_{kd}\hat{\epsilon}_{kd}$, $\hat{\epsilon} = \{\hat{\epsilon}_{kd}\}$
 - ▶ giving a K by D stochastic weight matrix: $\hat{W}(\theta, \hat{\epsilon})$
 - ▶ write $f^{\theta, \hat{\epsilon}}(x) = \hat{W}(\theta, \hat{\epsilon})^T \phi(x)$
 - ▶ giving

$$\hat{L}(\theta, \{\hat{\epsilon}_n\}) = -\frac{1}{2\sigma^2} \sum_{x_n, y_n} \|y_n - f^{\theta, \hat{\epsilon}_n}(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2s^2} \|M\|_2^2 \dots$$

- ▶ Model for regression (D outputs) / classification
- ▶ ELBO $L(\theta) = \int q_\theta(W) \log p(Y|X, W) dW - \text{KL}(q, \text{prior})$
- ▶ log likelihood eg

$$\log p(Y|X, W) = -\frac{1}{2\sigma^2} \sum \|y_n - f^W(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2$$
- ▶ approx post eg $q_\theta(w_{kd}) = N(w_{kd}; m_{kd}, \sigma_{kd}^2)$
 - ▶ $\text{KL}(q, \text{prior}) = \sum_{kd} 1/2(s^{-2}\sigma_{kd}^2 + s^{-2}m_{kd}^2 - 1 + \log(s^2/\sigma_{kd}^2))$
- ▶ MC integration:
 - ▶ sample $\hat{\epsilon}_{kd} \sim \mathcal{N}(0, 1)$ and write $\hat{w}_{kd} = m_{kd} + \sigma_{kd}\hat{\epsilon}_{kd}$, $\hat{\epsilon} = \{\hat{\epsilon}_{kd}\}$
 - ▶ giving a K by D stochastic weight matrix: $\hat{W}(\theta, \hat{\epsilon})$
 - ▶ write $f^{\theta, \hat{\epsilon}}(x) = \hat{W}(\theta, \hat{\epsilon})^T \phi(x)$
 - ▶ giving

$$\hat{L}(\theta, \{\hat{\epsilon}_n\}) = -\frac{1}{2\sigma^2} \sum_{x_n, y_n} \|y_n - f^{\theta, \hat{\epsilon}_n}(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2s^2} \|M\|_2^2 \dots$$

- ▶ until now we only did inference over W (last layer weights)
- ▶ because doing inference on preceding layers was too challenging (intractable / non-conjugate)
- ▶ but with our new techniques we can **easily** extend to W, b of all layers in model (denoted ω)
- ▶ these models (where all layers have dists over) are known as Bayesian neural networks (BNNs)
 - ▶ X of dim N by Q (and Y of dim N by D)
 - ▶ W^1 of dim Q by K , b^1 dim K
 - ▶ W^2 of dim K by D , b^2 dim D
 - ▶ ϕ elem-wise non-linearity
 - ▶ $\omega = \{W^1, W^2, b^1, b^2\}$
 - ▶ $f^\omega(x) = \phi(x^T W^1 + b^1) W^2 + b^2$
 - note:** could be a deep net with thousands of layers
- ▶ Long history (Hopfield [1987] \rightarrow LeCun [1991] \rightarrow MacKay [1992] \rightarrow Hinton [1993] \rightarrow Neal [1995] \rightarrow Barber and Bishop [1998])

- ▶ until now we only did inference over W (last layer weights)
- ▶ because doing inference on preceding layers was too challenging (intractable / non-conjugate)
- ▶ but with our new techniques we can **easily** extend to W, b of all layers in model (denoted ω)
- ▶ these models (where all layers have dists over) are known as Bayesian neural networks (BNNs)
 - ▶ X of dim N by Q (and Y of dim N by D)
 - ▶ W^1 of dim Q by K , b^1 dim K
 - ▶ W^2 of dim K by D , b^2 dim D
 - ▶ ϕ elem-wise non-linearity
 - ▶ $\omega = \{W^1, W^2, b^1, b^2\}$
 - ▶ $f^\omega(x) = \phi(x^T W^1 + b^1) W^2 + b^2$
note: could be a deep net with thousands of layers
- ▶ Long history (Hopfield [1987] \rightarrow LeCun [1991] \rightarrow MacKay [1992] \rightarrow Hinton [1993] \rightarrow Neal [1995] \rightarrow Barber and Bishop [1998])

- ▶ until now we only did inference over W (last layer weights)
- ▶ because doing inference on preceding layers was too challenging (intractable / non-conjugate)
- ▶ but with our new techniques we can **easily** extend to W, b of all layers in model (denoted ω)
- ▶ these models (where all layers have dists over) are known as Bayesian neural networks (BNNs)
 - ▶ X of dim N by Q (and Y of dim N by D)
 - ▶ W^1 of dim Q by K , b^1 dim K
 - ▶ W^2 of dim K by D , b^2 dim D
 - ▶ ϕ elem-wise non-linearity
 - ▶ $\omega = \{W^1, W^2, b^1, b^2\}$
 - ▶ $f^\omega(x) = \phi(x^T W^1 + b^1) W^2 + b^2$
note: could be a deep net with thousands of layers
- ▶ Long history (Hopfield [1987] \rightarrow LeCun [1991] \rightarrow MacKay [1992] \rightarrow Hinton [1993] \rightarrow Neal [1995] \rightarrow Barber and Bishop [1998])

BNNs

- ▶ model
 - ▶ as before, but swap W^1, W^2, f^ω instead of W and f^W
- ▶ approx inference
 - ▶ log likelihood – same
 - ▶ approx post – Gaussians w means $\{m_{qk}^1, m_{kd}^2\}$ and stds $\{\sigma_{qk}^1, \sigma_{kd}^2\}$
 - KL to prior $KL(q(W^1, W^2), p) = KL(q(W^1), p) + KL(q(W^2), p)$
 - ▶ ELBO

$$\hat{W}_{qk}^1 = m_{qk}^1 + \sigma_{qk}^1 \hat{\epsilon}_{qk}^1$$

$$\hat{W}_{kd}^2 = m_{kd}^2 + \sigma_{kd}^2 \hat{\epsilon}_{kd}^2$$

with $\hat{\epsilon}_{qk}^1, \hat{\epsilon}_{kd}^2 \sim \mathcal{N}(0, 1)$ and $\hat{\epsilon} = \{\hat{\epsilon}_{qk}^1, \hat{\epsilon}_{kd}^2\}$

- ▶ swap $f^{\theta, \epsilon}(x) = \hat{W}(\theta, \hat{\epsilon})^T \phi(x)$ with

$$f^{\theta, \epsilon}(x) = \phi(x^T \hat{W}^1(\theta, \hat{\epsilon}) + b^1) \hat{W}^2(\theta, \hat{\epsilon}) + b^2$$

and plug into $\hat{L}(\theta, \{\hat{\epsilon}_n\})$.

- ▶ Proposed as *MDL* from compression literature [Graves, 2011], in Bayesian modelling known as *mean-field variational inference (MFVI)*; Also referred to as *Bayes by Backprop* [Blundell, 2015].

BNNs

- ▶ model
 - ▶ as before, but swap W^1, W^2, f^ω instead of W and f^W
- ▶ approx inference
 - ▶ log likelihood – same
 - ▶ approx post – Gaussians w means $\{m_{qk}^1, m_{kd}^2\}$ and stds $\{\sigma_{qk}^1, \sigma_{kd}^2\}$

$$\text{KL to prior } \text{KL}(q(W^1, W^2), p) = \text{KL}(q(W^1), p) + \text{KL}(q(W^2), p)$$

- ▶ ELBO

$$\hat{W}_{qk}^1 = m_{qk}^1 + \sigma_{qk}^1 \hat{\epsilon}_{qk}^1$$

$$\hat{W}_{kd}^2 = m_{kd}^2 + \sigma_{kd}^2 \hat{\epsilon}_{kd}^2$$

with $\hat{\epsilon}_{qk}^1, \hat{\epsilon}_{kd}^2 \sim \mathcal{N}(0, 1)$ and $\hat{\epsilon} = \{\hat{\epsilon}_{qk}^1, \hat{\epsilon}_{kd}^2\}$

- ▶ swap $f^{\theta, \hat{\epsilon}}(x) = \hat{W}(\theta, \hat{\epsilon})^T \phi(x)$ with

$$f^{\theta, \hat{\epsilon}}(x) = \phi(x^T \hat{W}^1(\theta, \hat{\epsilon}) + b^1) \hat{W}^2(\theta, \hat{\epsilon}) + b^2$$

and plug into $\hat{L}(\theta, \{\hat{\epsilon}_n\})$.

- ▶ Proposed as *MDL* from compression literature [Graves, 2011], in Bayesian modelling known as *mean-field variational inference (MFVI)*; Also referred to as *Bayes by Backprop* [Blundell, 2015].

BNNs

- ▶ model
 - ▶ as before, but swap W^1, W^2, f^ω instead of W and f^W
- ▶ approx inference
 - ▶ log likelihood – same
 - ▶ approx post – Gaussians w means $\{m_{qk}^1, m_{kd}^2\}$ and stds $\{\sigma_{qk}^1, \sigma_{kd}^2\}$
 KL to prior $KL(q(W^1, W^2), p) = KL(q(W^1), p) + KL(q(W^2), p)$
 - ▶ ELBO

$$\hat{W}_{qk}^1 = m_{qk}^1 + \sigma_{qk}^1 \hat{\epsilon}_{qk}^1$$

$$\hat{W}_{kd}^2 = m_{kd}^2 + \sigma_{kd}^2 \hat{\epsilon}_{kd}^2$$

with $\hat{\epsilon}_{qk}^1, \hat{\epsilon}_{kd}^2 \sim \mathcal{N}(0, 1)$ and $\hat{\epsilon} = \{\hat{\epsilon}_{qk}^1, \hat{\epsilon}_{kd}^2\}$

- ▶ swap $f^{\theta, \hat{\epsilon}}(x) = \hat{W}(\theta, \hat{\epsilon})^T \phi(x)$ with

$$f^{\theta, \hat{\epsilon}}(x) = \phi(x^T \hat{W}^1(\theta, \hat{\epsilon}) + b^1) \hat{W}^2(\theta, \hat{\epsilon}) + b^2$$

and plug into $\hat{L}(\theta, \{\hat{\epsilon}_n\})$.

- ▶ Proposed as *MDL* from compression literature [Graves, 2011], in Bayesian modelling known as *mean-field variational inference (MFVI)*; Also referred to as *Bayes by Backprop* [Blundell, 2015].

BNNs

- ▶ model
 - ▶ as before, but swap W^1, W^2, f^ω instead of W and f^W
- ▶ approx inference
 - ▶ log likelihood – same
 - ▶ approx post – Gaussians w means $\{m_{qk}^1, m_{kd}^2\}$ and stds $\{\sigma_{qk}^1, \sigma_{kd}^2\}$
 KL to prior $KL(q(W^1, W^2), p) = KL(q(W^1), p) + KL(q(W^2), p)$
 - ▶ ELBO

$$\hat{W}_{qk}^1 = m_{qk}^1 + \sigma_{qk}^1 \hat{\epsilon}_{qk}^1$$

$$\hat{W}_{kd}^2 = m_{kd}^2 + \sigma_{kd}^2 \hat{\epsilon}_{kd}^2$$

with $\hat{\epsilon}_{qk}^1, \hat{\epsilon}_{kd}^2 \sim \mathcal{N}(0, 1)$ and $\hat{\epsilon} = \{\hat{\epsilon}_{qk}^1, \hat{\epsilon}_{kd}^2\}$

- ▶ swap $f^{\theta, \hat{\epsilon}}(x) = \hat{W}(\theta, \hat{\epsilon})^T \phi(x)$ with

$$f^{\theta, \hat{\epsilon}}(x) = \phi(x^T \hat{W}^1(\theta, \hat{\epsilon}) + b^1) \hat{W}^2(\theta, \hat{\epsilon}) + b^2$$

and plug into $\hat{L}(\theta, \{\hat{\epsilon}_n\})$.

- ▶ Proposed as *MDL* from compression literature [Graves, 2011], in Bayesian modelling known as *mean-field variational inference (MFVI)*; Also referred to as *Bayes by Backprop* [Blundell, 2015].

BNNs

- ▶ model
 - ▶ as before, but swap W^1, W^2, f^ω instead of W and f^W
- ▶ approx inference
 - ▶ log likelihood – same
 - ▶ approx post – Gaussians w means $\{m_{qk}^1, m_{kd}^2\}$ and stds $\{\sigma_{qk}^1, \sigma_{kd}^2\}$
 - KL to prior $KL(q(W^1, W^2), p) = KL(q(W^1), p) + KL(q(W^2), p)$
 - ▶ ELBO

$$\hat{W}_{qk}^1 = m_{qk}^1 + \sigma_{qk}^1 \hat{\epsilon}_{qk}^1$$

$$\hat{W}_{kd}^2 = m_{kd}^2 + \sigma_{kd}^2 \hat{\epsilon}_{kd}^2$$

with $\hat{\epsilon}_{qk}^1, \hat{\epsilon}_{kd}^2 \sim \mathcal{N}(0, 1)$ and $\hat{\epsilon} = \{\hat{\epsilon}_{qk}^1, \hat{\epsilon}_{kd}^2\}$

- ▶ swap $f^{\theta, \hat{\epsilon}}(x) = \hat{W}(\theta, \hat{\epsilon})^T \phi(x)$ with

$$f^{\theta, \hat{\epsilon}}(x) = \phi(x^T \hat{W}^1(\theta, \hat{\epsilon}) + b^1) \hat{W}^2(\theta, \hat{\epsilon}) + b^2$$

and plug into $\hat{L}(\theta, \{\hat{\epsilon}_n\})$.

- ▶ Proposed as *MDL* from compression literature [Graves, 2011], in Bayesian modelling known as *mean-field variational inference (MFVI)*; Also referred to as *Bayes by Backprop* [Blundell, 2015].

Issue with above...

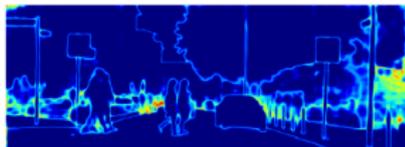
- ▶ when we use large models we usually use 10s-100s of millions of params – models as big as can fit on GPU
- ▶ when using Gaussian approx we need at least two params for each NN weight
- ▶ doubling num of params... so having to reduce model size by 2!
- ▶ can we scale the ideas above to very large models?



(a) Input Image



(b) Semantic Segmentation

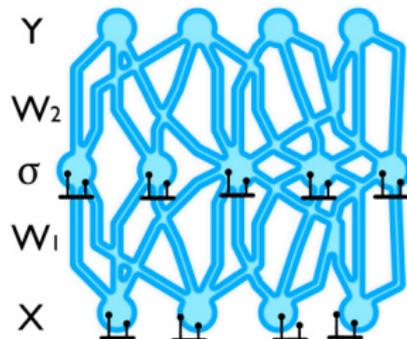


(c) Epistemic Uncertainty

Inference in Very Large Deep Models

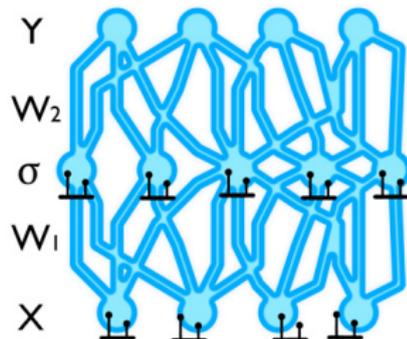
Stochastic regularisation

- ▶ lots of techniques in deep learning inject noise into large models to help with regularisation
- ▶ eg dropout (but lots of others which mostly work the same)
 - ▶ at training time, randomly set network units to zero with prob p (Bern)
 - ▶ call this “stochastic forward pass”
 - ▶ at test time multiply each unit by $1/(1 - p)$ and do not drop
 - ▶ call this “deterministic forward pass”
 - ▶ noise is added to units (feature space)
 - ▶ implemented in every deep learning framework (from TF/PyTorch to TensorRT for embedded devices)



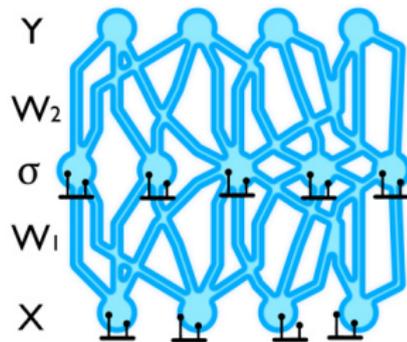
Stochastic regularisation

- ▶ lots of techniques in deep learning inject noise into large models to help with regularisation
- ▶ eg dropout (but lots of others which mostly work the same)
 - ▶ at training time, randomly set network units to zero with prob p (Bern)
 - ▶ call this “stochastic forward pass”
 - ▶ at test time multiply each unit by $1/(1 - p)$ and do not drop
 - ▶ call this “deterministic forward pass”
 - ▶ noise is added to units (feature space)
 - ▶ implemented in every deep learning framework (from TF/PyTorch to TensorRT for embedded devices)



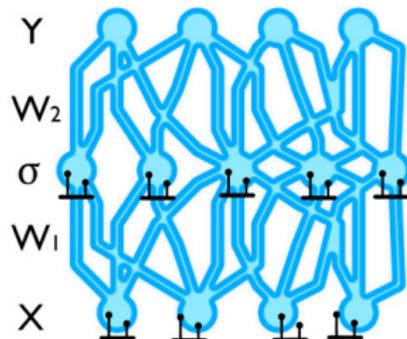
Stochastic regularisation

- ▶ lots of techniques in deep learning inject noise into large models to help with regularisation
- ▶ eg dropout (but lots of others which mostly work the same)
 - ▶ at training time, randomly set network units to zero with prob p (Bern)
 - ▶ call this “stochastic forward pass”
 - ▶ at test time multiply each unit by $1/(1 - p)$ and do not drop
 - ▶ call this “deterministic forward pass”
- ▶ noise is added to units (feature space)
- ▶ implemented in *every* deep learning framework (from TF/PyTorch to TensorRT for embedded devices)



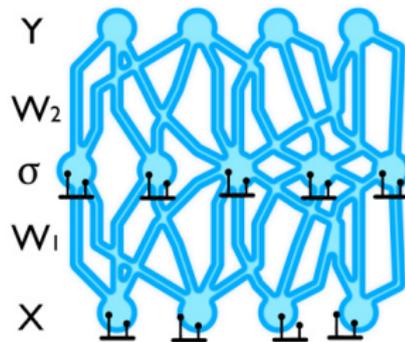
Stochastic regularisation

- ▶ lots of techniques in deep learning inject noise into large models to help with regularisation
- ▶ eg dropout (but lots of others which mostly work the same)
 - ▶ at training time, randomly set network units to zero with prob p (Bern)
 - ▶ call this “stochastic forward pass”
 - ▶ at test time multiply each unit by $1/(1 - p)$ and do not drop
 - ▶ call this “deterministic forward pass”
 - ▶ noise is added to units (feature space)
 - ▶ implemented in *every* deep learning framework (from TF/PyTorch to TensorRT for embedded devices)



Stochastic regularisation

- ▶ lots of techniques in deep learning inject noise into large models to help with regularisation
- ▶ eg dropout (but lots of others which mostly work the same)
 - ▶ at training time, randomly set network units to zero with prob p (Bern)
 - ▶ call this “stochastic forward pass”
 - ▶ at test time multiply each unit by $1/(1 - p)$ and do not drop
 - ▶ call this “deterministic forward pass”
 - ▶ noise is added to units (feature space)
 - ▶ implemented in *every* deep learning framework (from TF/PyTorch to TensorRT for embedded devices)



$$X : \mathbb{R}^1$$

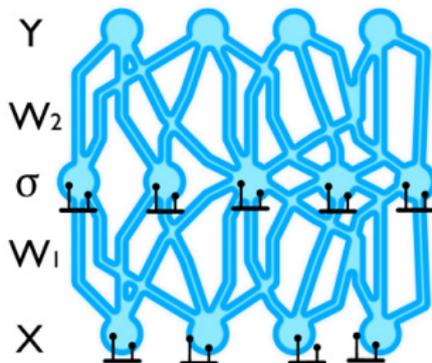
$$\hat{X} := x^T \hat{E}^1, \quad [\hat{E}^1]_{ij} \sim \text{Bern}(p^1), \quad \text{a } 2 \text{ by } 2 \text{ matrix (zero off diag)}$$

$$h := \phi(x^T M^1 + b^1)$$

$$\hat{h} := h \hat{E}^2, \quad [\hat{E}^2]_{kk} \sim \text{Bern}(p^2)$$

$$\hat{y} := \hat{h} M^2 + b^2$$

$$\hat{J} = \frac{1}{N} \sum_n \|y_n - \hat{y}_n\|_2^2 + \lambda_1 \|M^1\|_2^2 + \lambda_2 \|M^2\|_2^2$$



- ▶ Can transform noise to param (weight) space

$$\begin{aligned}\hat{y} &= \left(\phi \left[(x\hat{\epsilon}^1)M^1 + b^1 \right] \hat{\epsilon}^2 \right) M^2 + b^2 \\ &= \phi \left[x(\hat{\epsilon}^1 M^1) + b^1 \right] (\hat{\epsilon}^2 M^2) + b^2\end{aligned}$$

writing $\hat{W}^1 := \hat{\epsilon}^1 M^1$ and $\hat{W}^2 := \hat{\epsilon}^2 M^2$ gives

$$\hat{y} = \phi(x\hat{W}^1 + b^1)\hat{W}^2 + b^2 = f^{\hat{\omega}}(x)$$

with $\hat{\omega} = \{\hat{W}^1, \hat{W}^2\}$

- ▶ so at training time dropout samples weights matrices... looks v familiar!
- ▶ let's see if we can make this connection more formal
- ▶ develop approx inf in BNNs with a new approx dist...

- ▶ Can transform noise to param (weight) space

$$\begin{aligned}\hat{y} &= \left(\phi \left[(x\hat{\epsilon}^1)M^1 + b^1 \right] \hat{\epsilon}^2 \right) M^2 + b^2 \\ &= \phi \left[x(\hat{\epsilon}^1 M^1) + b^1 \right] (\hat{\epsilon}^2 M^2) + b^2\end{aligned}$$

writing $\hat{W}^1 := \hat{\epsilon}^1 M^1$ and $\hat{W}^2 := \hat{\epsilon}^2 M^2$ gives

$$\hat{y} = \phi(x\hat{W}^1 + b^1)\hat{W}^2 + b^2 = f^{\hat{\omega}}(x)$$

with $\hat{\omega} = \{\hat{W}^1, \hat{W}^2\}$

- ▶ so at training time dropout samples weights matrices... looks v familiar!
- ▶ let's see if we can make this connection more formal
- ▶ develop approx inf in BNNs with a new approx dist...

- ▶ model

- ▶ prior - same
- ▶ lik - same

- ▶ approx inference

- ▶ approx dist $q_\theta(W^1) = M^1 \epsilon$ with $\epsilon_{gg} = \text{Bernoulli}(p^1)$ and zero otherwise, and $\theta = \{M^1, p^1, M^2, p^2\}$

$$\text{KL}(q, p) \approx \frac{1-p^1}{2s^2} \|M^1\|_2^2 - QH(p^1) + \frac{1-p^2}{2s^2} \|M^2\|_2^2 - KH(p^2) + \text{const}$$

- ▶ ELBO

$$\hat{L}(\theta, \{\hat{\epsilon}_n\}) = -\frac{1}{2\sigma^2} \sum_{x_n, y_n} \|y_n - f^{\theta, \hat{\epsilon}_n}(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2 - \frac{1-p^1}{2s^2} \|M^1\|_2^2$$

with $f^{\theta, \hat{\epsilon}_n}(x_n)$ a dropout stochastic forward pass

- ▶ can rewrite as min obj (multiply by $-2\sigma^2/N$)

$$J = \frac{1}{N} \sum_n \|y_n - \hat{y}_n\|_2^2 + \lambda^1 \|M^1\|_2^2 + \lambda^2 \|M^2\|_2^2 + \text{const}$$

with $\hat{y}_n = f^{\theta, \hat{\epsilon}_n}(x_n)$ and defining $\lambda^1 = \sigma^2 \frac{1-p^1}{s^2 N}$..

- ▶ model
 - ▶ prior - same
 - ▶ lik - same
- ▶ approx inference
 - ▶ approx dist $q_\theta(W^1) = M^1 \epsilon$ with $\epsilon_{qq} = \text{Bernoulli}(p^1)$ and zero otherwise, and $\theta = \{M^1, p^1, M^2, p^2\}$

$$\text{KL}(q, p) \approx \frac{1-p^1}{2s^2} \|M^1\|_2^2 - \text{QH}(p^1) + \frac{1-p^2}{2s^2} \|M^2\|_2^2 - \text{KH}(p^2) + \text{const}$$

- ▶ ELBO

$$\hat{L}(\theta, \{\hat{\epsilon}_n\}) = -\frac{1}{2\sigma^2} \sum_{x_n, y_n} \|y_n - f^{\theta, \hat{\epsilon}_n}(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2 - \frac{1-p^1}{2s^2} \|M^1\|_2^2$$

with $f^{\theta, \hat{\epsilon}_n}(x_n)$ a dropout stochastic forward pass

- ▶ can rewrite as min obj (multiply by $-2\sigma^2/N$)

$$J = \frac{1}{N} \sum_n \|y_n - \hat{y}_n\|_2^2 + \lambda^1 \|M^1\|_2^2 + \lambda^2 \|M^2\|_2^2 + \text{const}$$

with $\hat{y}_n = f^{\theta, \hat{\epsilon}_n}(x_n)$ and defining $\lambda^1 = \sigma^2 \frac{1-p^1}{s^2 N}$..

- ▶ model
 - ▶ prior - same
 - ▶ lik - same
- ▶ approx inference
 - ▶ approx dist $q_\theta(W^1) = M^1 \epsilon$ with $\epsilon_{qq} = \text{Bernoulli}(p^1)$ and zero otherwise, and $\theta = \{M^1, p^1, M^2, p^2\}$

$$\text{KL}(q, p) \approx \frac{1-p^1}{2s^2} \|M^1\|_2^2 - \text{QH}(p^1) + \frac{1-p^2}{2s^2} \|M^2\|_2^2 - \text{KH}(p^2) + \text{const}$$

- ▶ ELBO

$$\hat{L}(\theta, \{\hat{\epsilon}_n\}) = -\frac{1}{2\sigma^2} \sum_{x_n, y_n} \|y_n - f^{\theta, \hat{\epsilon}_n}(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2 - \frac{1-p^1}{2s^2} \|M^1\|_2^2$$

with $f^{\theta, \hat{\epsilon}_n}(x_n)$ a dropout stochastic forward pass

- ▶ can rewrite as min obj (multiply by $-2\sigma^2/N$)

$$J = \frac{1}{N} \sum_n \|y_n - \hat{y}_n\|_2^2 + \lambda^1 \|M^1\|_2^2 + \lambda^2 \|M^2\|_2^2 + \text{const}$$

with $\hat{y}_n = f^{\theta, \hat{\epsilon}_n}(x_n)$ and defining $\lambda^1 = \sigma^2 \frac{1-p^1}{s^2 N} \dots$

- ▶ model
 - ▶ prior - same
 - ▶ lik - same
- ▶ approx inference
 - ▶ approx dist $q_\theta(W^1) = M^1 \epsilon$ with $\epsilon_{qq} = \text{Bernoulli}(p^1)$ and zero otherwise, and $\theta = \{M^1, p^1, M^2, p^2\}$

$$\text{KL}(q, p) \approx \frac{1-p^1}{2s^2} \|M^1\|_2^2 - QH(p^1) + \frac{1-p^2}{2s^2} \|M^2\|_2^2 - KH(p^2) + \text{const}$$

- ▶ ELBO

$$\hat{L}(\theta, \{\hat{\epsilon}_n\}) = -\frac{1}{2\sigma^2} \sum_{x_n, y_n} \|y_n - f^{\theta, \hat{\epsilon}_n}(x_n)\|_2^2 - \frac{N}{2} \log 2\pi\sigma^2 - \frac{1-p^1}{2s^2} \|M^1\|_2^2$$

with $f^{\theta, \hat{\epsilon}_n}(x_n)$ a dropout stochastic forward pass

- ▶ can rewrite as min obj (multiply by $-2\sigma^2/N$)

$$J = \frac{1}{N} \sum_n \|y_n - \hat{y}_n\|_2^2 + \lambda^1 \|M^1\|_2^2 + \lambda^2 \|M^2\|_2^2 + \text{const}$$

with $\hat{y}_n = f^{\theta, \hat{\epsilon}_n}(x_n)$ and defining $\lambda^1 = \sigma^2 \frac{1-p^1}{s^2 N}$..

- ▶ This is the standard dropout objective
- ▶ ie any standard NN in which you use dropout, you can view as a BNN
- ▶ note: need to tune p as a variational param
 - ▶ can't diff wrt p (used in Bern in obj; can't use reparam trick..)
 - ▶ but when you do grid search over p on a validation set, use $\hat{L}(\theta, \{\hat{\epsilon}_n\})$ to select p which max ELBO or validation log predictive
 - ▶ Can also use continuous relaxation for dropout (see Concrete Dropout, 2017)
- ▶ Example:

- ▶ This is the standard dropout objective
- ▶ ie any standard NN in which you use dropout, you can view as a BNN
- ▶ note: need to tune p as a variational param
 - ▶ can't diff wrt p (used in Bern in obj; can't use reparam trick..)
 - ▶ but when you do grid search over p on a validation set, use $\hat{L}(\theta, \{\hat{\epsilon}_n\})$ to select p which max ELBO or validation log predictive
 - ▶ Can also use continuous relaxation for dropout (see Concrete Dropout, 2017)
- ▶ Example:

- ▶ This is the standard dropout objective
- ▶ ie any standard NN in which you use dropout, you can view as a BNN
- ▶ note: need to tune p as a variational param
 - ▶ can't diff wrt p (used in Bern in obj; can't use reparam trick..)
 - ▶ but when you do grid search over p on a validation set, use $\hat{L}(\theta, \{\hat{\epsilon}_n\})$ to select p which max ELBO or validation log predictive
 - ▶ Can also use continuous relaxation for dropout (see Concrete Dropout, 2017)
- ▶ Example:

- ▶ This is the standard dropout objective
- ▶ ie any standard NN in which you use dropout, you can view as a BNN
- ▶ note: need to tune p as a variational param
 - ▶ can't diff wrt p (used in Bern in obj; can't use reparam trick..)
 - ▶ but when you do grid search over p on a validation set, use $\hat{L}(\theta, \{\hat{\epsilon}_n\})$ to select p which max ELBO or validation log predictive
 - ▶ Can also use continuous relaxation for dropout (see Concrete Dropout, 2017)
- ▶ Example:

Dropout uncertainty example

Define model and train on data `x_train`, `y_train`:

```
1 from tensorflow.keras.layers import Input, Dense, Dropout
2 from tf.keras.regularizers import l2
3
4 reg = sigma**2 * (1-p) / (s**2 * N)
5
6 inputs = Input(shape=(512,))
7 x = Dense(1024, activation="relu",
8           kernel_regularizer=l2(reg))(inputs)
9 x = Dropout(p)(x, training=True)
10 x = Dense(1024, activation="relu",
11           kernel_regularizer=l2(reg))(x)
12 x = Dropout(p)(x, training=True)
13 outputs = Dense(1, kernel_regularizer=l2(reg))(x)
14
15 model = tf.keras.Model(inputs, outputs)
16 model.compile(loss="mean_squared_error",
17              optimizer="adam")
18 model.fit(x_train, y_train)
```

Using MC estimators can estimate epistemic uncertainty in BNNs almost trivially...

- ▶ predictive mean



$$E_{p(y^*|x^*, \mathcal{D})}[y^*] \approx \frac{1}{T} \sum_t f^{\hat{\omega}_t}(x)$$

with $\hat{\omega}_t \sim q_{\theta}(\omega)$.

ie, average multiple stochastic forward passes

- ▶ predictive variance

- ▶ again, collect some stochastic forward passes...

$$\begin{aligned} \text{Var}_{p(y^*|x^*, \mathcal{D})}[y^*] &= E_{p(y^*|x^*, \mathcal{D})}[(y^*)^2] - E_{p(y^*|x^*, \mathcal{D})}[y^*]^2 \\ &\approx \sigma^2 + \frac{1}{T} \sum_t f^{\hat{\omega}_t}(x)^2 - \left(\frac{1}{T} \sum_t f^{\hat{\omega}_t}(x) \right)^2 \end{aligned}$$

Using MC estimators can estimate epistemic uncertainty in BNNs almost trivially...

► predictive mean



$$E_{p(y^*|x^*, \mathcal{D})}[y^*] \approx \frac{1}{T} \sum_t f^{\hat{\omega}_t}(x)$$

with $\hat{\omega}_t \sim q_{\theta}(\omega)$.

ie, average multiple stochastic forward passes

► predictive variance

► again, collect some stochastic forward passes...

$$\begin{aligned} \text{Var}_{p(y^*|x^*, \mathcal{D})}[y^*] &= E_{p(y^*|x^*, \mathcal{D})}[(y^*)^2] - E_{p(y^*|x^*, \mathcal{D})}[y^*]^2 \\ &\approx \sigma^2 + \frac{1}{T} \sum_t f^{\hat{\omega}_t}(x)^2 - \left(\frac{1}{T} \sum_t f^{\hat{\omega}_t}(x) \right)^2 \end{aligned}$$

Predictive mean:

$$\begin{aligned}
 & E_{p(y^* | x^*, D)} \{y^*\} \\
 &= \int p(y^* | x^*, D) y^* dy^* \\
 &= \int \left(\int p(y^* | x^*, \omega) p(\omega | D) d\omega \right) y^* dy^* \\
 &= \int \left(\underbrace{\int y^* p(y^* | x^*, \omega) dy^*}_{=f^w(x, \sigma^2)} \right) p(\omega | D) d\omega \\
 &= \int f^w(x) \underbrace{p(\omega | D)}_{\approx q(\omega)} d\omega
 \end{aligned}$$

$$\text{VI} \quad \hat{\mu} \approx \int f^w(x) q(\omega) d\omega$$

$$\text{Mc} \quad \hat{\mu} \approx \frac{1}{T} \sum_t f^{\tilde{w}_t}(x), \quad \tilde{w}_t \sim q(\omega)$$

Predictive Variance:

$$E_{p(y^* | x^*, D)} \{y^{*2}\}$$

$$= \int \left(\int y^{*2} p(y^* | x^*, w) dy^* \right) p(w | D) dw$$

$E_{p(y^* | x^*, D)} \{y^{*2}\} = \text{Var}_{p(y^* | x^*, D)} \{y^*\} + E_{p(y^* | x^*, D)} \{y^*\}^2$

$$\stackrel{\text{lik}}{=} \int (\sigma^2 + f^w(x)^2) p(w | D) dw$$

$$\stackrel{\text{VI}}{\approx} \sigma^2 + \int f^w(x)^2 q(w) dw$$

$$\stackrel{\text{MC}}{\approx} \sigma^2 + \frac{1}{T} \sum_t f^{\tilde{w}_t}(x)^2$$

$$\Rightarrow \text{Var}_{p(y^* | x^*, D)} \{y^*\} = E_{p(y^* | x^*, D)} \{y^{*2}\} - E_{p(y^* | x^*, D)} \{y^*\}^2$$

$$\approx \sigma^2 + \frac{1}{T} \sum_t f^{\tilde{w}_t}(x)^2 - \left(\frac{1}{T} \sum_t f^{\tilde{w}_t}(x) \right)^2$$

Do stochastic forward passes on `x_test`:

```
1 | num_MC_samples = 100
2 | MC_samples = [model.predict(x_test)
3 |                 for _ in range(num_MC_samples)]
```

Predictive mean

```
1 | np.mean(MC_samples, axis=0)
```

Predictive variance

```
1 | sigma**2 + np.var(MC_samples, axis=0)
```

A useful tool for debugging

- ▶ sample from weights $\omega \sim p(\omega|\mathcal{D}) =$ function sample $f^\omega(\cdot)$
- ▶ evaluate over interval $[-10, 10]$
- ▶ eg:
 - ▶ sample ω and def $f^\omega(\cdot)$
 - ▶ for each x_i in $\{-10, -9.95, -9.9, \dots, 9.9, 9.95, 10\}$
 - ▶ evaluate $y_i = f^\omega(x_i)$ and plot (x_i, y_i)
- ▶ note: if using dropout inference, use same dropout mask for all inputs x
- ▶ Visualisation: bd1101.ml/vis

A useful tool for debugging

- ▶ sample from weights $\omega \sim p(\omega|\mathcal{D}) =$ function sample $f^\omega(\cdot)$
- ▶ evaluate over interval $[-10, 10]$
- ▶ eg:
 - ▶ sample ω and def $f^\omega(\cdot)$
 - ▶ for each x_i in $\{-10, -9.95, -9.9, \dots, 9.9, 9.95, 10\}$
 - ▶ evaluate $y_i = f^\omega(x_i)$ and plot (x_i, y_i)
- ▶ note: if using dropout inference, use same dropout mask for all inputs x
- ▶ Visualisation: bd1101.ml/vis

A useful tool for debugging

- ▶ sample from weights $\omega \sim p(\omega|\mathcal{D}) =$ function sample $f^\omega(\cdot)$
- ▶ evaluate over interval $[-10, 10]$
- ▶ eg:
 - ▶ sample ω and def $f^\omega(\cdot)$
 - ▶ for each x_i in $\{-10, -9.95, -9.9, \dots, 9.9, 9.95, 10\}$
 - ▶ evaluate $y_i = f^\omega(x_i)$ and plot (x_i, y_i)
- ▶ note: if using dropout inference, use same dropout mask for all inputs x
- ▶ Visualisation: bd1101.ml/vis

Real-world Applications of Model Uncertainty

- ▶ we use machine learning to aid experts working in laborious fields
- ▶ automate small parts of the expert's work
 - ▶ eg melanoma (cancer) diagnosis based on lesion images
- ▶ but deep learning often requires large amounts of labelled data
 - ▶ increases with the complexity of problem
 - ▶ complexity of the input data
 - ▶ eg image inputs require large models
 - ▶ hundreds of gigabytes in ImageNet
- ▶ sometimes can't afford to label huge data...
 - ▶ eg automating lesion image analysis
 - ▶ would require expert to spend expensive time annotating large number of lesion images (for every cancer type of interest)
- ▶ instead, could use *active learning*



- ▶ we use machine learning to aid experts working in laborious fields
- ▶ automate small parts of the expert's work
 - ▶ eg melanoma (cancer) diagnosis based on lesion images
- ▶ but deep learning often requires large amounts of labelled data
 - ▶ increases with the complexity of problem
 - ▶ complexity of the input data
 - ▶ eg image inputs require large models
 - ▶ hundreds of gigabytes in ImageNet
- ▶ sometimes can't afford to label huge data...
 - ▶ eg automating lesion image analysis
 - ▶ would require expert to spend expensive time annotating large number of lesion images (for every cancer type of interest)
- ▶ instead, could use *active learning*



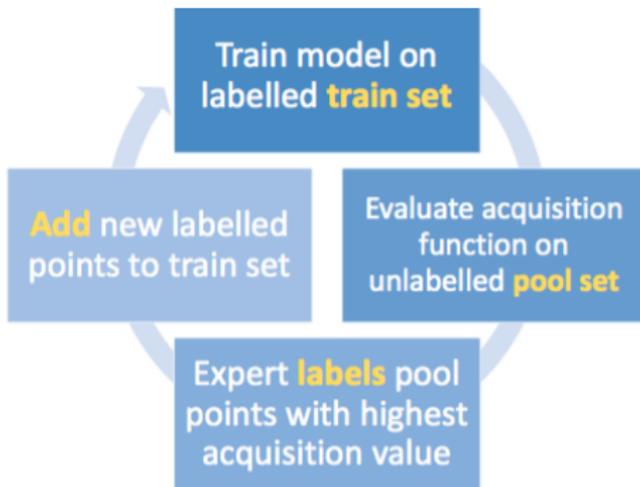
- ▶ we use machine learning to aid experts working in laborious fields
- ▶ automate small parts of the expert's work
 - ▶ eg melanoma (cancer) diagnosis based on lesion images
- ▶ but deep learning often requires large amounts of labelled data
 - ▶ increases with the complexity of problem
 - ▶ complexity of the input data
 - ▶ eg image inputs require large models
 - ▶ hundreds of gigabytes in ImageNet
- ▶ sometimes can't afford to label huge data...
 - ▶ eg automating lesion image analysis
 - ▶ would require expert to spend expensive time annotating large number of lesion images (for every cancer type of interest)
- ▶ instead, could use *active learning*

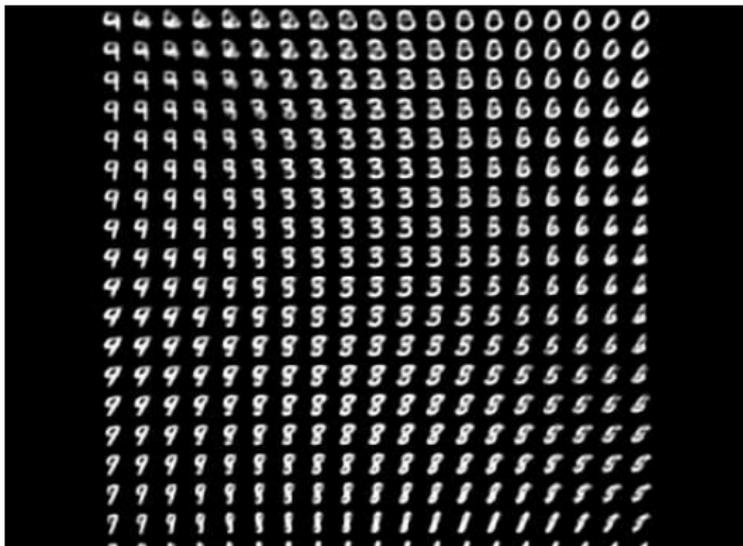


- ▶ we use machine learning to aid experts working in laborious fields
- ▶ automate small parts of the expert's work
 - ▶ eg melanoma (cancer) diagnosis based on lesion images
- ▶ but deep learning often requires large amounts of labelled data
 - ▶ increases with the complexity of problem
 - ▶ complexity of the input data
 - ▶ eg image inputs require large models
 - ▶ hundreds of gigabytes in ImageNet
- ▶ sometimes can't afford to label huge data...
 - ▶ eg automating lesion image analysis
 - ▶ would require expert to spend expensive time annotating large number of lesion images (for every cancer type of interest)
- ▶ instead, could use *active learning*



- ▶ active learning
 - ▶ *agent* chooses which unlabelled data is most informative
 - ▶ asks external “oracle” (eg human annotator) for a label only for that
 - ▶ acquisition function: ranks points based on their potential informativeness
 - ▶ eg, epistemic uncertainty





```
1 ...  
2 model.compile(loss="categorical_crossentropy",  
3               optimizer="adam")  
4 model.fit(x_train, y_train)  
5 MC_samples = [model.predict(x_test) for _ in range(20)]
```

Need uncertainty for classification...

```

1 | ...
2 | model.compile(loss="categorical_crossentropy",
3 |               optimizer="adam")
4 | model.fit(x_train, y_train)
5 | MC_samples = [model.predict(x_test) for _ in range(20)]
  
```

Predictive entropy

$$p(y^* = c | x^*, \mathcal{D}) \approx \frac{1}{T} \sum_t p^{\hat{W}_t}(x^*)_c$$

$$H_{p(y^* | x^*, \mathcal{D})}[y^*] = - \sum_{y^*=c} p(y^* = c | x^*, \mathcal{D}) \log p(y^* = c | x^*, \mathcal{D})$$

```

1 | expected_p = np.mean(MC_samples, axis=0)
2 | predictive_entropy = -np.sum(expected_p *
3 |                               np.log(expected_p), axis=-1)
  
```

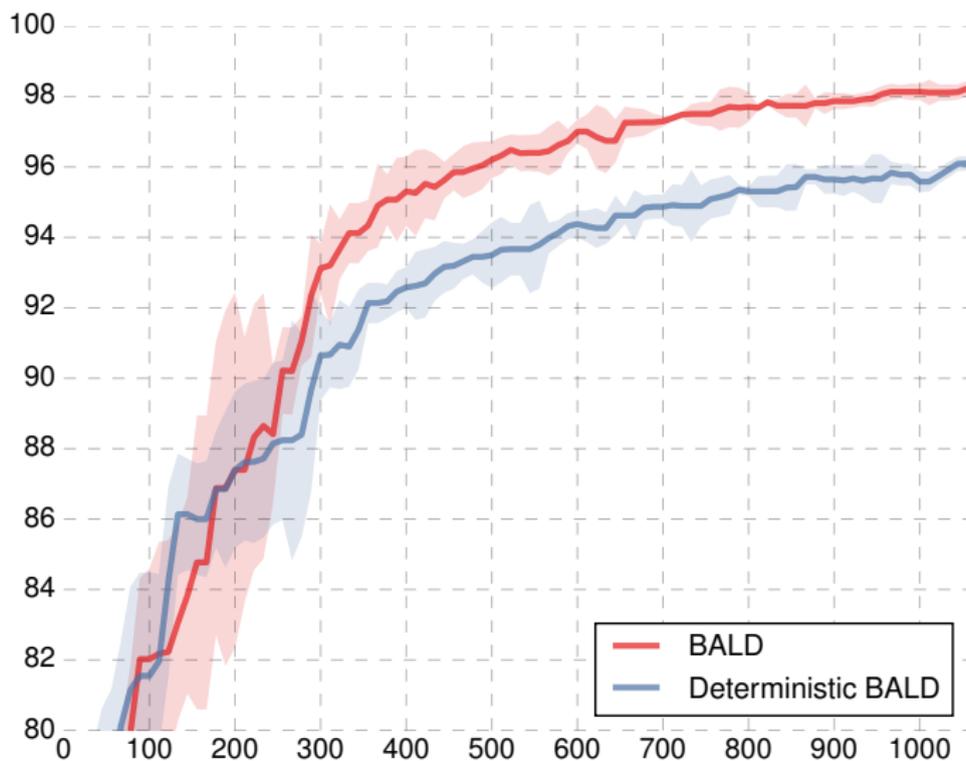
```
1 | ...
2 | model.compile(loss="categorical_crossentropy",
3 |               optimizer="adam")
4 | model.fit(x_train, y_train)
5 | MC_samples = [model.predict(x_test) for _ in range(20)]
```

Mutual information (epistemic uncertainty)

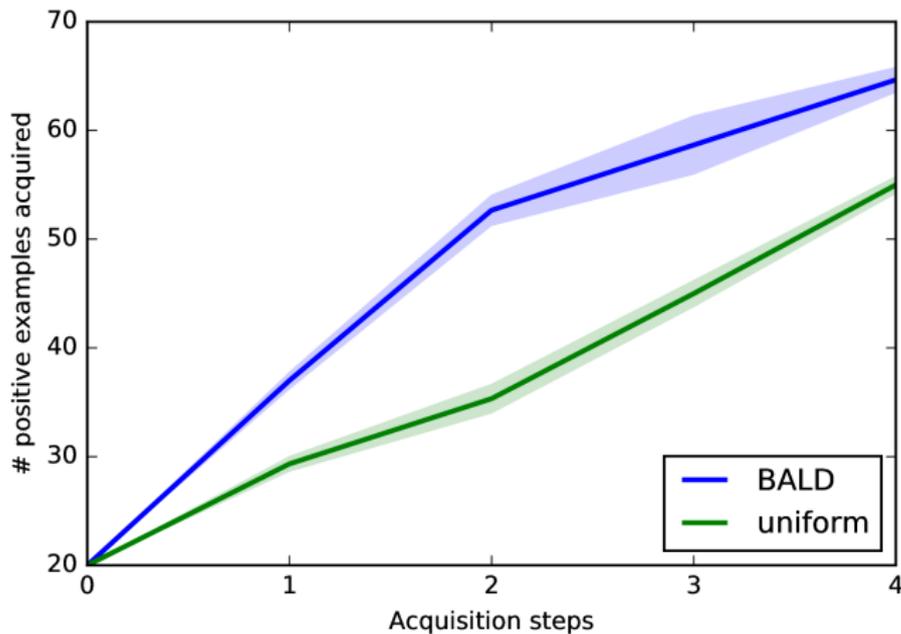
$$MI(y^*, W|\mathcal{D}, x^*) = H_{p(y^*|x^*, \mathcal{D})}[y^*] - \frac{1}{T} \sum_{t, y^*=c} p^{\hat{W}_t(x^*)}_c \log p^{\hat{W}_t(x^*)}_c$$

```
1 | MC_entropy = np.sum(MC_samples * np.log(MC_samples),
2 |                    axis=-1)
3 | expected_entropy = -np.mean(MC_entropy, axis=0)
4 | mi = predictive_entropy - expected_entropy
```

MNIST with only 1,000 images (instead of 60,000)



Melanoma diagnosis with 300 images



acquired positive examples vs. acquisition



- ▶ goal is to detect diabetes, and be able to tell when model is guessing at random
- ▶ used to pre-screen patients, send only patients with high uncertainty to expert



```
t = tf.layers.Conv2D(initial_conv_
                    kernel_regula
t = tf.layers.dropout(t, rate=drop
t = tf.layers.MaxPooling2D(pool_s

t = tf.layers.Conv2D(initial_conv_
                    kernel_regula
t = tf.layers.dropout(t, rate=drop
t = tf.layers.Conv2D(initial_conv_
                    padding="same
t = tf.layers.dropout(t, rate=drop
t = tf.layers.MaxPooling2D(pool_s
```

How to tell if uncertainty is good or bad?

- ▶ define a binary event: 'is diabetes?'; group test set inputs by **prediction** 'yes'/'no' vs **label** 'yes'/'no'
- ▶ each corresponds to one of TP, FP, FN, TN
- ▶ TPR and FPR are rates of TP and FP
 - ▶ $\text{TPR} = \text{sensitivity} = \text{recall} = \text{TP} / (\text{TP} + \text{FN}) = 1 - \text{FNR}$
 - ▶ $\text{TNR} = \text{specificity} = \text{TN} / (\text{TN} + \text{FP}) = 1 - \text{FPR}$
 - ▶ $\text{FPR} = \text{FP} / (\text{TN} + \text{FP}) = 1 - \text{specificity}$
- ▶ want TPR to be high, FPR to be low
- ▶ usually given reqs what's the worst we're allowed to perform in order to deploy system
- ▶ eg TPR=0.7 and FPR=0.1

How to tell if uncertainty is good or bad?

- ▶ define a binary event: 'is diabetes?'; group test set inputs by **prediction** 'yes'/'no' vs **label** 'yes'/'no'
- ▶ each corresponds to one of TP, FP, FN, TN
- ▶ TPR and FPR are rates of TP and FP
 - ▶ $\text{TPR} = \text{sensitivity} = \text{recall} = \text{TP} / (\text{TP} + \text{FN}) = 1 - \text{FNR}$
 - ▶ $\text{TNR} = \text{specificity} = \text{TN} / (\text{TN} + \text{FP}) = 1 - \text{FPR}$
 - ▶ $\text{FPR} = \text{FP} / (\text{TN} + \text{FP}) = 1 - \text{specificity}$
- ▶ want TPR to be high, FPR to be low
- ▶ usually given reqs what's the worst we're allowed to perform in order to deploy system
- ▶ eg TPR=0.7 and FPR=0.1

How to tell if uncertainty is good or bad?

- ▶ define a binary event: 'is diabetes?'; group test set inputs by **prediction** 'yes'/'no' vs **label** 'yes'/'no'
- ▶ each corresponds to one of TP, FP, FN, TN
- ▶ TPR and FPR are rates of TP and FP
 - ▶ $\text{TPR} = \text{sensitivity} = \text{recall} = \text{TP} / (\text{TP} + \text{FN}) = 1 - \text{FNR}$
 - ▶ $\text{TNR} = \text{specificity} = \text{TN} / (\text{TN} + \text{FP}) = 1 - \text{FPR}$
 - ▶ $\text{FPR} = \text{FP} / (\text{TN} + \text{FP}) = 1 - \text{specificity}$
- ▶ want TPR to be high, FPR to be low
- ▶ usually given reqs what's the worst we're allowed to perform in order to deploy system
- ▶ eg TPR=0.7 and FPR=0.1

How to tell if uncertainty is good or bad?

- ▶ define a binary event: 'is diabetes?'; group test set inputs by **prediction** 'yes'/'no' vs **label** 'yes'/'no'
- ▶ each corresponds to one of TP, FP, FN, TN
- ▶ TPR and FPR are rates of TP and FP
 - ▶ $\text{TPR} = \text{sensitivity} = \text{recall} = \text{TP} / (\text{TP} + \text{FN}) = 1 - \text{FNR}$
 - ▶ $\text{TNR} = \text{specificity} = \text{TN} / (\text{TN} + \text{FP}) = 1 - \text{FPR}$
 - ▶ $\text{FPR} = \text{FP} / (\text{TN} + \text{FP}) = 1 - \text{specificity}$
- ▶ want TPR to be high, FPR to be low
- ▶ usually given reqs what's the worst we're allowed to perform in order to deploy system
- ▶ eg TPR=0.7 and FPR=0.1

How to tell if uncertainty is good or bad?

- ▶ define a binary event: 'is diabetes?'; group test set inputs by **prediction** 'yes'/'no' vs **label** 'yes'/'no'
- ▶ each corresponds to one of TP, FP, FN, TN
- ▶ TPR and FPR are rates of TP and FP
 - ▶ $\text{TPR} = \text{sensitivity} = \text{recall} = \text{TP} / (\text{TP} + \text{FN}) = 1 - \text{FNR}$
 - ▶ $\text{TNR} = \text{specificity} = \text{TN} / (\text{TN} + \text{FP}) = 1 - \text{FPR}$
 - ▶ $\text{FPR} = \text{FP} / (\text{TN} + \text{FP}) = 1 - \text{specificity}$
- ▶ want TPR to be high, FPR to be low
- ▶ usually given reqs what's the worst we're allowed to perform in order to deploy system
- ▶ eg TPR=0.7 and FPR=0.1

- ▶ model outputs a predictive prob $p(y|x, D)$; how do we get a recommendation 'yes'/'no'?
 - ▶ easiest is to take argmax
 - ▶ but what if model outputs 0.51? is this a 'yes'?
- ▶ def a threshold t
- ▶ if predictive prob is higher than t then say 'yes' otherwise say 'no'
 - ▶ for $t = 0$ says 'yes' to all, ie $FN=TN=0$, and model has $TPR=1, FPR=1$
 - ▶ for $t = 1$ says 'no' to all, ie $TP=FP=0$ and model has $TPR=0, FPR=0$
- ▶ each threshold t gives us a pair (FPR, TPR)
- ▶ scatter points for all t (or some discrete steps t)
- ▶ this is an ROC plot

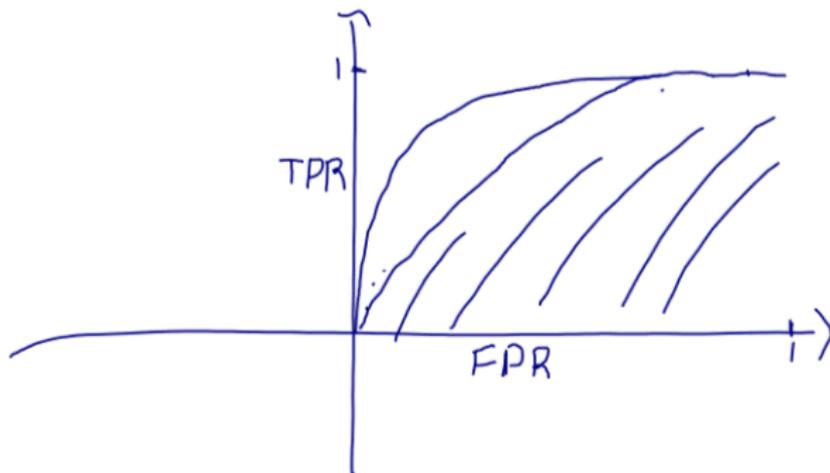
- ▶ model outputs a predictive prob $p(y|x, D)$; how do we get a recommendation 'yes'/'no'?
 - ▶ easiest is to take argmax
 - ▶ but what if model outputs 0.51? is this a 'yes'?
- ▶ def a threshold t
- ▶ if predictive prob is higher than t then say 'yes' otherwise say 'no'
 - ▶ for $t = 0$ says 'yes' to all, ie $FN=TN=0$, and model has $TPR=1, FPR=1$
 - ▶ for $t = 1$ says 'no' to all, ie $TP=FP=0$ and model has $TPR=0, FPR=0$
- ▶ each threshold t gives us a pair (FPR, TPR)
- ▶ scatter points for all t (or some discrete steps t)
- ▶ this is an ROC plot

- ▶ model outputs a predictive prob $p(y|x, D)$; how do we get a recommendation 'yes'/'no'?
 - ▶ easiest is to take argmax
 - ▶ but what if model outputs 0.51? is this a 'yes'?
- ▶ def a threshold t
- ▶ if predictive prob is higher than t then say 'yes' otherwise say 'no'
 - ▶ for $t = 0$ says 'yes' to all, ie $FN=TN=0$, and model has $TPR=1, FPR=1$
 - ▶ for $t = 1$ says 'no' to all, ie $TP=FP=0$ and model has $TPR=0, FPR=0$
- ▶ each threshold t gives us a pair (FPR, TPR)
- ▶ scatter points for all t (or some discrete steps t)
- ▶ this is an ROC plot

- ▶ model outputs a predictive prob $p(y|x, D)$; how do we get a recommendation 'yes'/'no'?
 - ▶ easiest is to take argmax
 - ▶ but what if model outputs 0.51? is this a 'yes'?
- ▶ def a threshold t
- ▶ if predictive prob is higher than t then say 'yes' otherwise say 'no'
 - ▶ for $t = 0$ says 'yes' to all, ie $FN=TN=0$, and model has $TPR=1, FPR=1$
 - ▶ for $t = 1$ says 'no' to all, ie $TP=FP=0$ and model has $TPR=0, FPR=0$
- ▶ each threshold t gives us a pair (FPR, TPR)
- ▶ scatter points for all t (or some discrete steps t)
- ▶ this is an ROC plot

- ▶ model outputs a predictive prob $p(y|x, D)$; how do we get a recommendation 'yes'/'no'?
 - ▶ easiest is to take argmax
 - ▶ but what if model outputs 0.51? is this a 'yes'?
- ▶ def a threshold t
- ▶ if predictive prob is higher than t then say 'yes' otherwise say 'no'
 - ▶ for $t = 0$ says 'yes' to all, ie $FN=TN=0$, and model has $TPR=1, FPR=1$
 - ▶ for $t = 1$ says 'no' to all, ie $TP=FP=0$ and model has $TPR=0, FPR=0$
- ▶ each threshold t gives us a pair (FPR, TPR)
- ▶ scatter points for all t (or some discrete steps t)
- ▶ this is an ROC plot

- ▶ model outputs a predictive prob $p(y|x, D)$; how do we get a recommendation 'yes'/'no'?
 - ▶ easiest is to take argmax
 - ▶ but what if model outputs 0.51? is this a 'yes'?
- ▶ def a threshold t
- ▶ if predictive prob is higher than t then say 'yes' otherwise say 'no'
 - ▶ for $t = 0$ says 'yes' to all, ie $FN=TN=0$, and model has $TPR=1, FPR=1$
 - ▶ for $t = 1$ says 'no' to all, ie $TP=FP=0$ and model has $TPR=0, FPR=0$
- ▶ each threshold t gives us a pair (FPR, TPR)
- ▶ scatter points for all t (or some discrete steps t)
- ▶ this is an ROC plot



- ▶ ROC shows tradeoff between TPR and FPR
- ▶ each point on the plot corresponds to a choice of t which will give that tradeoff
- ▶ aim: find a model which gives highest Area Under Curve (AUC)
 - ▶ allows for better tradeoffs generally
 - ▶ but not always
- ▶ how can we improve AUC? one solution:
 - ▶ identify patients for which you are guessing at random (uncertain)
 - ▶ select 10% patients you are most uncertain about and remove from test set (send to expert)
 - ▶ plot ROC for remaining 90% test set patients
 - ▶ if uncertainty correlates to patients you were mistaken on, ROC should improve (higher AUC)

- ▶ ROC shows tradeoff between TPR and FPR
- ▶ each point on the plot corresponds to a choice of t which will give that tradeoff
- ▶ aim: find a model which gives highest Area Under Curve (AUC)
 - ▶ allows for better tradeoffs generally
 - ▶ but not always
- ▶ how can we improve AUC? one solution:
 - ▶ identify patients for which you are guessing at random (uncertain)
 - ▶ select 10% patients you are most uncertain about and remove from test set (send to expert)
 - ▶ plot ROC for remaining 90% test set patients
 - ▶ if uncertainty correlates to patients you were mistaken on, ROC should improve (higher AUC)

- ▶ ROC shows tradeoff between TPR and FPR
- ▶ each point on the plot corresponds to a choice of t which will give that tradeoff
- ▶ aim: find a model which gives highest Area Under Curve (AUC)
 - ▶ allows for better tradeoffs generally
 - ▶ but not always
- ▶ how can we improve AUC? one solution:
 - ▶ identify patients for which you are guessing at random (uncertain)
 - ▶ select 10% patients you are most uncertain about and remove from test set (send to expert)
 - ▶ plot ROC for remaining 90% test set patients
 - ▶ if uncertainty correlates to patients you were mistaken on, ROC should improve (higher AUC)

- ▶ ROC shows tradeoff between TPR and FPR
- ▶ each point on the plot corresponds to a choice of t which will give that tradeoff
- ▶ aim: find a model which gives highest Area Under Curve (AUC)
 - ▶ allows for better tradeoffs generally
 - ▶ but not always
- ▶ how can we improve AUC? one solution:
 - ▶ identify patients for which you are guessing at random (uncertain)
 - ▶ select 10% patients you are most uncertain about and remove from test set (send to expert)
 - ▶ plot ROC for remaining 90% test set patients
 - ▶ if uncertainty correlates to patients you were mistaken on, ROC should improve (higher AUC)

Send patients to expert diagnosis if model is uncertain

- ▶ use some uncertainty metric to refuse to diagnose a patient if model is uncertain
- ▶ what uncertainty measure?
 - ▶ MI would be high for far away points but will keep ambiguous points in test set
 - ▶ (points for which expert annotation in dataset was noisy)
 - ▶ expected entropy would be high for both far away inputs (entropy \geq MI) and ambiguous inputs
- ▶ → use expected entropy
- ▶ can we improve tradeoff by sending a small number of patients to an expert in a real-world system?

Send patients to expert diagnosis if model is uncertain

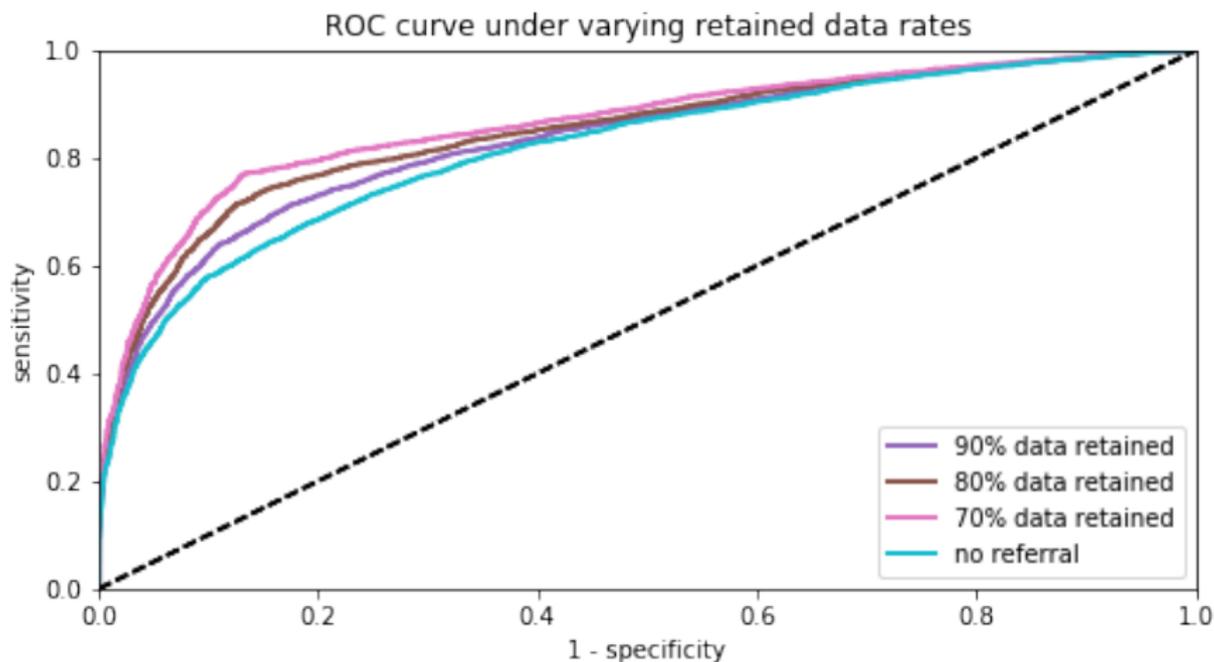
- ▶ use some uncertainty metric to refuse to diagnose a patient if model is uncertain
- ▶ what uncertainty measure?
 - ▶ MI would be high for far away points but will keep ambiguous points in test set
 - ▶ (points for which expert annotation in dataset was noisy)
 - ▶ expected entropy would be high for both far away inputs (entropy \geq MI) and ambiguous inputs
- ▶ → use expected entropy
- ▶ can we improve tradeoff by sending a small number of patients to an expert in a real-world system?

Send patients to expert diagnosis if model is uncertain

- ▶ use some uncertainty metric to refuse to diagnose a patient if model is uncertain
- ▶ what uncertainty measure?
 - ▶ MI would be high for far away points but will keep ambiguous points in test set
 - ▶ (points for which expert annotation in dataset was noisy)
 - ▶ expected entropy would be high for both far away inputs (entropy \geq MI) and ambiguous inputs
- ▶ → use expected entropy
- ▶ can we improve tradeoff by sending a small number of patients to an expert in a real-world system?

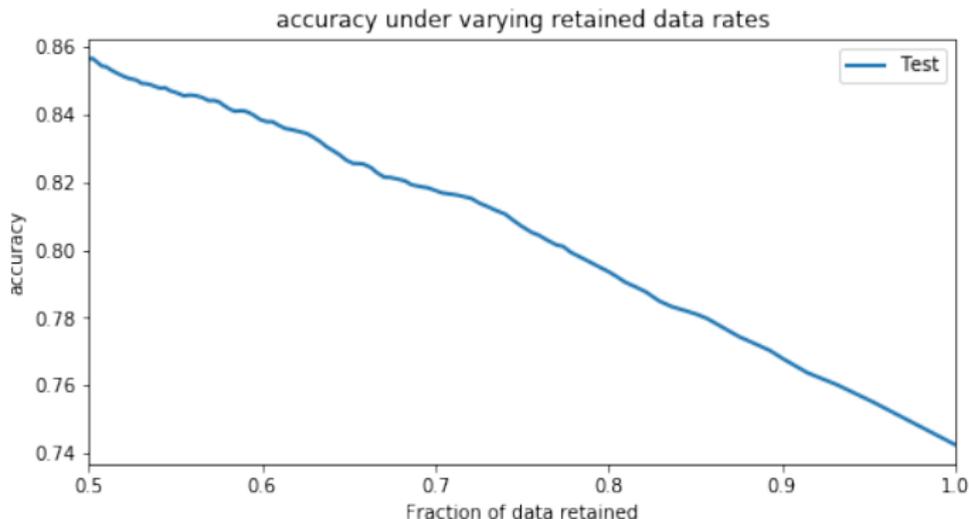
Send patients to expert diagnosis if model is uncertain

- ▶ use some uncertainty metric to refuse to diagnose a patient if model is uncertain
- ▶ what uncertainty measure?
 - ▶ MI would be high for far away points but will keep ambiguous points in test set
 - ▶ (points for which expert annotation in dataset was noisy)
 - ▶ expected entropy would be high for both far away inputs (entropy \geq MI) and ambiguous inputs
- ▶ → use expected entropy
- ▶ can we improve tradeoff by sending a small number of patients to an expert in a real-world system?



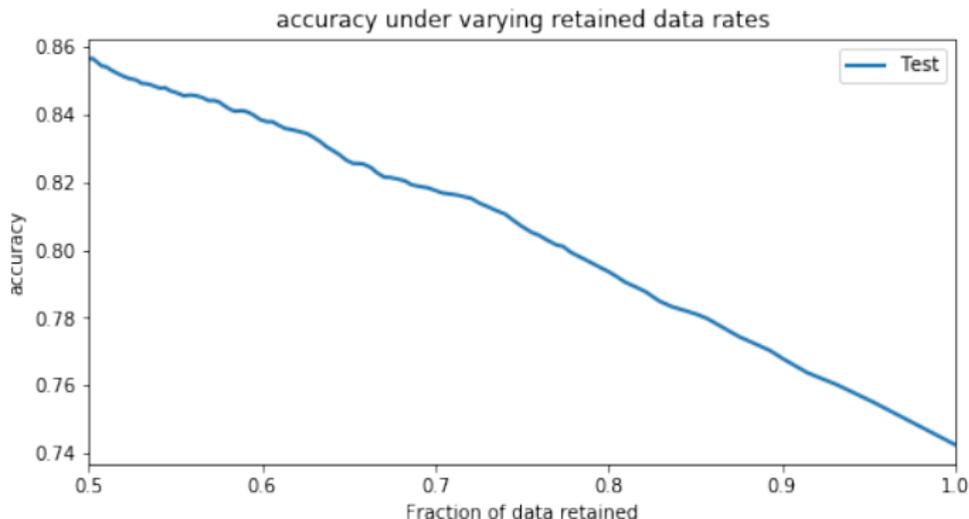
Another measure of uncertainty performance

- ▶ plot accuracy as a function of % retained data, as sending more and more patients to an expert
- ▶ 100% retain data = original accuracy on full dataset
- ▶ 10% retain data = accuracy after removing 90% patients with highest uncertainty



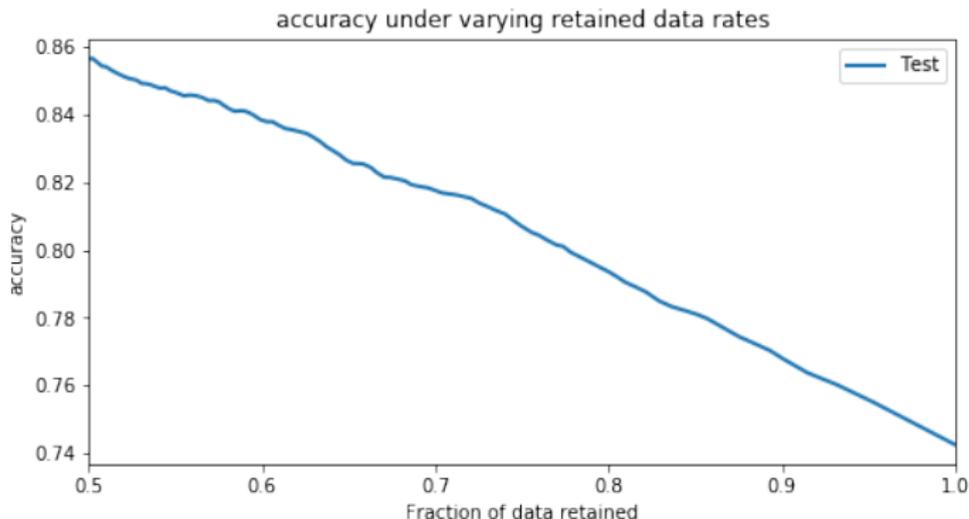
Another measure of uncertainty performance

- ▶ plot accuracy as a function of % retained data, as sending more and more patients to an expert
- ▶ 100% retain data = original accuracy on full dataset
- ▶ 10% retain data = accuracy after removing 90% patients with highest uncertainty



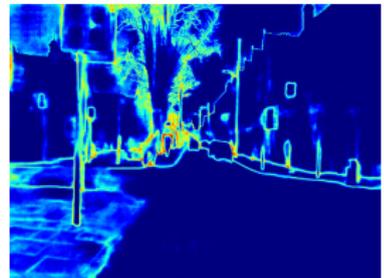
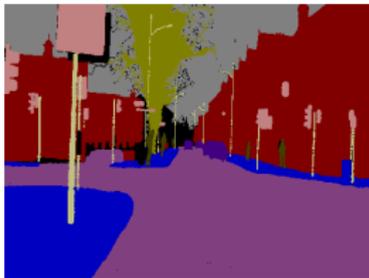
Another measure of uncertainty performance

- ▶ plot accuracy as a function of % retained data, as sending more and more patients to an expert
- ▶ 100% retain data = original accuracy on full dataset
- ▶ 10% retain data = accuracy after removing 90% patients with highest uncertainty



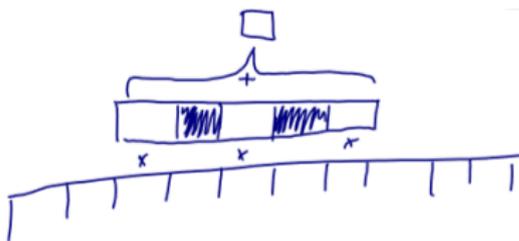
We'll be looking at semantic segmentation

- ▶ input: image in RGB space
- ▶ output: image in semantic space
- ▶ each pixel is mapped to semantic class (eg road, sky, car, pedestrian) based on its context (near by pixels)



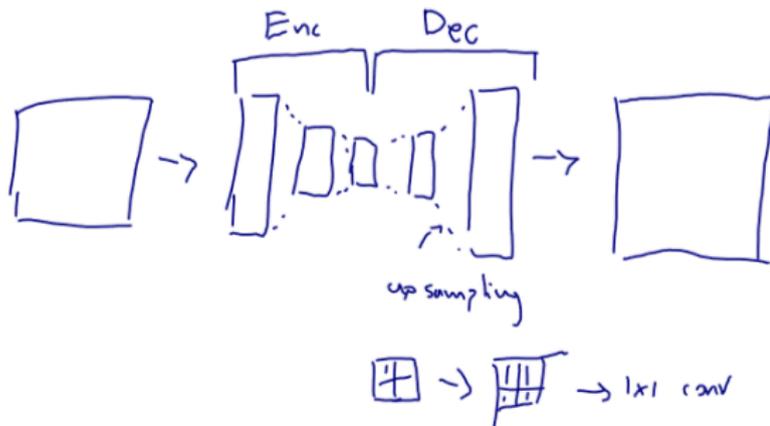
- ▶ one of the SOTA NNs for semantic segmentation is DeepLab
- ▶ uses atrous (dilated) convolutions (has 'holes')
 - ▶ widen field of view over the input feature maps without increasing parameters or pooling
- ▶ uses encoder-decoder architectures
 - ▶ upsampling replicates pixels then applies eg 1x1 conv which doesn't reduce dim
- ▶ can be applied to any base network ('backbone') as long as it is fully convolutional (ie no fully connected layers)

- ▶ one of the SOTA NNs for semantic segmentation is DeepLab
- ▶ uses atrous (dilated) convolutions (has 'holes')
- ▶ widen field of view over the input feature maps without increasing parameters or pooling



- ▶ uses encoder-decoder architectures
 - ▶ upsampling replicates pixels then applies eg 1x1 conv which doesn't reduce dim
- ▶ can be applied to any base network ('backbone') as long as it is fully convolutional (ie no fully connected layers)

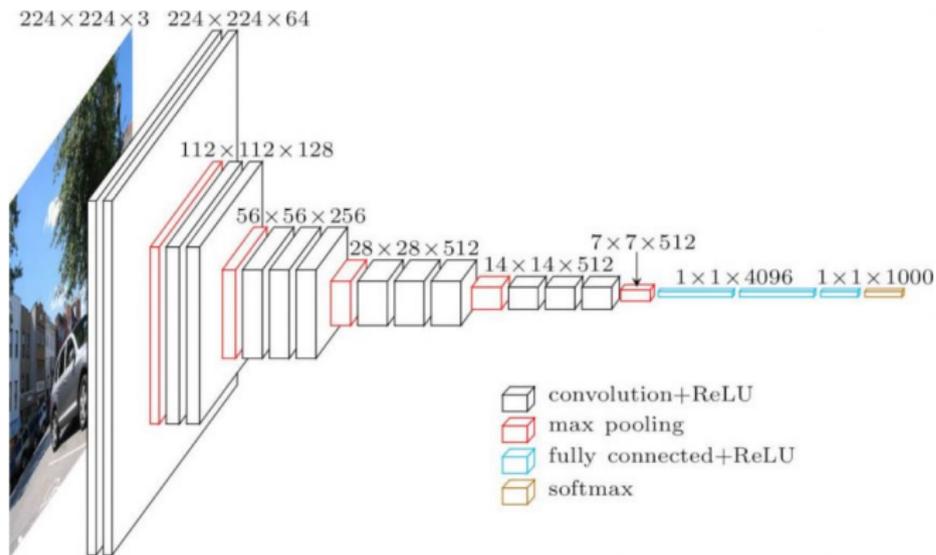
- ▶ one of the SOTA NNs for semantic segmentation is DeepLab
- ▶ uses atrous (dilated) convolutions (has 'holes')
 - ▶ widen field of view over the input feature maps without increasing parameters or pooling
- ▶ uses encoder-decoder architectures
 - ▶ upsampling replicates pixels then applies eg 1×1 conv which doesn't reduce dim



- ▶ one of the SOTA NNs for semantic segmentation is DeepLab
- ▶ uses atrous (dilated) convolutions (has 'holes')
 - ▶ widen field of view over the input feature maps without increasing parameters or pooling
- ▶ uses encoder-decoder architectures
 - ▶ upsampling replicates pixels then applies eg 1x1 conv which doesn't reduce dim
- ▶ can be applied to any base network ('backbone') as long as it is fully convolutional (ie no fully connected layers)

► Popular deep CNNs backbones

- VGG-16
- ResNet101
- Xception



neurohive.io/

► ResNet

- ▶ Popular deep CNNs backbones
 - ▶ VGG-16
 - ▶ ResNet101
 - ▶ Xception

- ▶ ResNet
 - ▶ layer def
 - ▶ solves the issue of “diminishing gradient” in deep nets (bounding eigenvalues from below)
 - ▶ can use hundreds of layers - seems to improve results the more layers you use

$$f^{w'}(x) = \overset{Q \times Q}{\downarrow} (w' + \mathbb{I}_Q) x = w'x + x$$

$$y(x) = w f^{w'}(f^{w^2}(f^{w^3}(x)))$$

$$= w(w' + \mathbb{I})(w^2 + \mathbb{I})(w^3 + \mathbb{I})x$$

$$\rightarrow = w w' w^2 w^3 x + w w^2 w^3 x + w w^3 x + \dots$$

- ▶ Popular deep CNNs backbones
 - ▶ VGG-16
 - ▶ ResNet101
 - ▶ Xception

- ▶ ResNet
 - ▶ layer def
 - ▶ solves the issue of “diminishing gradient” in deep nets (bounding eigenvalues from below)
 - ▶ can use hundreds of layers - seems to improve results the more layers you use

$$L = (y - y(x))^2$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y(x)} \frac{\partial y(x)}{\partial w}$$

$\underbrace{\qquad\qquad\qquad}_{-2(y-y(x))}$
 $\underbrace{\qquad\qquad\qquad}_{?}$

Deep Net:

$$y(x) = w f(w' (L^{w^2} (f^{w^3}(x)))) = w w' w^2 w^3 x$$

$$\frac{\partial y(x)}{\partial w^3} = w w' w^2 x$$

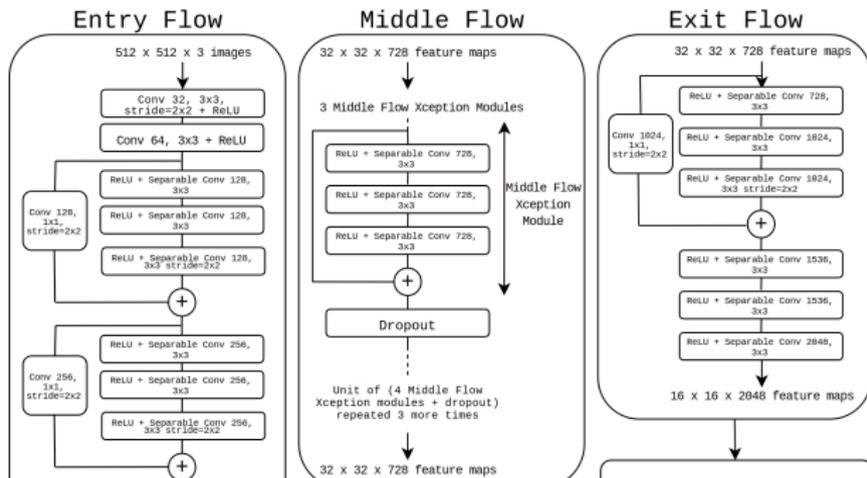
ResNet:

$$\frac{\partial y(x)}{\partial w^3} = w w' w^2 x + w w^2 x + w w' x + w x$$

$$= w(w w' w^2 + w^2 + w' + 1) x$$

We use Xception

- ▶ architecture has simplicity of VGG with multiple convolution layers stacked on top of one another
- ▶ Xception modules use skip connections similar to ResNet but between blocks
- ▶ works well empirically



- ▶ we have a classification problem with H by W softmax outputs (categorical variable for each pixel)
- ▶ model loss: sum of cross entropy (log likelihoods) for each pixel
- ▶ can use standard tools for uncertainty in classification (per pixel)
- ▶ and look at epistemic and aleatoric uncertainty maps

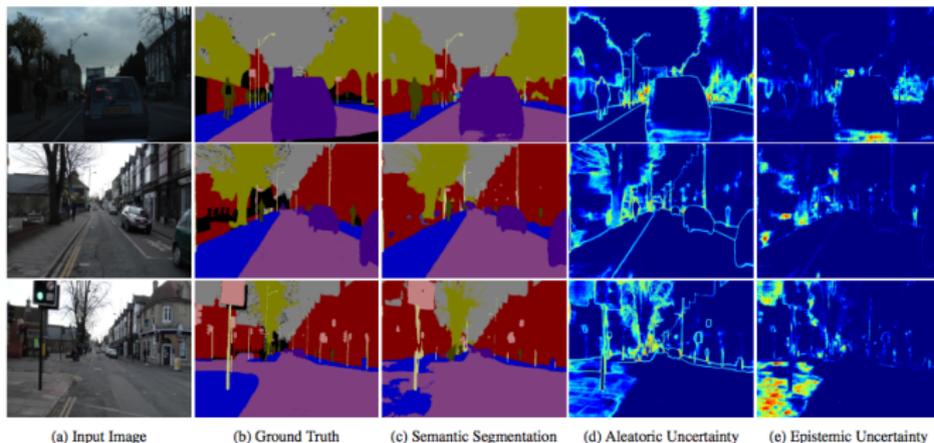


Figure 1. Qualitative results of our method on the CamVid semantic segmentation dataset, with three example input images. Our

www.nature.com/scientificreports

SCIENTIFIC REPORTS

OPEN

Leveraging uncertainty information from deep neural networks for disease detection

Christian Lebig¹, Vaneeda Allken¹, Murat Seçkin Ayhan¹, Philipp Berens^{1,2} & Siegfried Wahl^{1,3}

: 24 July 2017

: 1 December 2017

online: 19 December 2017

Deep learning (DL) has revolutionized the field of computer vision and image processing. In medical imaging, algorithmic solutions based on DL have been shown to achieve high performance on tasks that previously required medical experts. However, DL-based solutions for disease detection have been proposed without methods to quantify and control their uncertainty in a decision. In contrast, a physician knows whether she is uncertain about a case and will consult more experienced colleagues if

SC

Uncertainty-Aware Reinforcement Learning for Collision Avoidance

 Gregory Kahn*, Adam Villafior*, Vitchyr Pong*, Pieter Abbeel*[†], Sergey Levine*

*Berkeley AI Research (BAIR), University of California, Berkeley

[†]OpenAI

Abstract—Reinforcement learning can enable complex, adaptive behavior to be learned automatically for autonomous robotic platforms. However, practical deployment of reinforcement learning methods must contend with the fact that the training process itself can be unsafe for the robot. In this paper, we consider the specific case of a mobile robot learning to navigate an a priori unknown environment while avoiding collisions. In order to learn collision avoidance, the robot must experience collisions at training time. However, high-speed collisions, even at training time, could damage the robot. A successful learning method must therefore proceed cautiously, experiencing only low-speed collisions until it gains confidence. To this end, we present an uncertainty-aware model-based learning algorithm that estimates the probability of collision together with a statistical estimate of uncertainty. By formulating an uncertainty-dependent cost function, we show that the algorithm naturally chooses to proceed cautiously in unfamiliar environments, and increases the velocity of the robot in settings where it has high confidence. Our predictive model is based on bootstrapped neural networks

: 24 July 2017

: 1 December 2017

online: 19 Decem



Fig. 1: **Uncertainty-aware collision prediction model for collision avoidance:** A quadrotor and an RC car are tasked with navigating in an unknown environment. How should the robots navigate while avoiding collisions? We propose a model-based reinforcement learning approach in which the robot learns a collision prediction model by experiencing collisions at low speed which is unlikely to damage the vehicle. We formulate a velocity-dependent collision cost that uses collision prediction estimates and their associated uncertainties to enable the robot to only experience safe collisions during training while still approaching the desired task performance.

catastrophic) collisions during training. The robot can overcome this quandary by first experiencing gentle collisions in order

Uncertainty-Aware Reinforcement Learning for



This CVPR workshop paper is the Open Access version, provided by the Computer Vision Foundation. Except for this watermark, it is identical to the version available on IEEE Xplore.

SC

Semantic Segmentation of Small Objects and Modeling of Uncertainty in Urban Remote Sensing Images Using Deep Convolutional Neural Networks

O Abstract—Reinforcement learning (RL) has shown promising behavior to learn on complex environments. However, RL methods must learn on a specific case before they can be used on a new environment. In this paper, we propose a method to learn on a specific case before training time. We propose a method to learn on a specific case before training time. We propose a method to learn on a specific case before training time.

: 24 July 2017
: 1 December 2017
online: 19 December 2017

Michael Kampffmeyer*, Arnt-Børre Salberg† and Robert Jenssen*

*Machine Learning @ UiT Lab, UiT–The Arctic University of Norway

†Norwegian Computing Center

Abstract

We propose a deep Convolutional Neural Network (CNN) for land cover mapping in remote sensing images, with a focus on urban areas. In remote sensing, class imbalance represents often a problem for tasks like land cover

Remote sensing imagery is often characterized by complex data properties in the form of heterogeneity and class imbalance, as well as overlapping class-conditional distributions [6]. Together, these aspects constitute severe challenges for creating land cover maps or detecting and localizing objects, producing a high degree of uncertainty in ob-

SC

Uncertainty-Aware Reinforcement Learning for



THE CVF... Vision Foundation. Xplore.

Semantic
Ren

THE ASTROPHYSICAL JOURNAL LETTERS

certainty in Urban
ral Networks

: 24 July 2017
: 1 December 2017
l online: 19 Decem

O Abstract—Rein-
ive behavior to l
latforms. How
ng methods mus
self can be un
he specific case
riori unknown
o learn collision
t training time.
ime, could dan
ust therefore p
ollisions until it
ncertainty-awa
he probability
f uncertainty. I
unction, we shov
autiously in un
f the robot in
redictive mode

*We propose
(CNN) for land
with a focus on
balance represe.*

Uncertainties in Parameters Estimated with Neural Networks: Application to Strong Gravitational Lensing

Laurence Perreault Levasseur^{1,2} ,
Yashar D. Hezaveh^{1,2,3} , and
Risa H. Wechsler^{1,2} 

Published 2017 November 15 • © 2017. The
American Astronomical Society. All rights reserved.

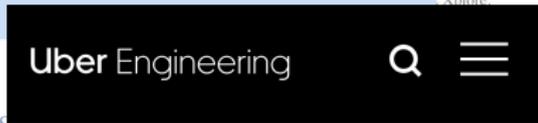
issen*
Norway

s often characterized by com-
rm of heterogeneity and class
pping class-conditional distri-
spects constitute severe chal-
: maps or detecting and local-
h degree of uncertainty in ob-

SC

Uncertainty-Aware Reinforcement Learning for

CvF



Abstract—Reinforcement learning (RL) is a powerful framework for learning behavior to interact with dynamic environments. However, existing RL methods must learn the specific case of the environment a priori, which is often unknown or difficult to learn at training time. In this paper, we propose uncertainty-aware RL, which learns the probability of collisions until the robot reaches a collision. We show that uncertainty-aware RL can be used to autonomously in unmodeled environments. We propose a method for learning the probability of collisions in an uncertainty-aware RL framework. We show that uncertainty-aware RL can be used to autonomously in unmodeled environments. We propose a method for learning the probability of collisions in an uncertainty-aware RL framework. We show that uncertainty-aware RL can be used to autonomously in unmodeled environments.

July 24, 2017

December 1, 2017

December 19, 2017

Semantic
Ren

THE AS

in Urban
rks

Uber Data

Uncertainty
Estimation
Neural Networks
Gravitational

Engineering Uncertainty Estimation in Neural Networks for Time Series Prediction at Uber

challenged by com-
plexity and class
distributional distri-
bution. We address
these challenges by
learning a model that
is robust to these
challenges and local-
izing uncertainty in ob-

We propose a method for learning the probability of collisions in an uncertainty-aware RL framework. We show that uncertainty-aware RL can be used to autonomously in unmodeled environments.

Laurenc
Yashar I
Risa H. I
Publisher
American

Lingxue Zhu and Nikolay Laptev

September 6,
2017



- ▶ use uncertainty in regression **correctly**
- ▶ perform predictions in simple probabilistic models **efficiently**
- ▶ use Bayesian modelling in complex ML models (eg classification)
- ▶ use uncertainty (both epistemic and aleatoric) in real world models
- ▶ extend VI **correctly** to complex models
 - ▶ try to extend to new likelihoods like Laplace
 - ▶ try to extend to multiple outputs: categorical and continuous outputs
- ▶ do deep learning with small amounts of data
 - ▶ do try this at home!
- ▶ evaluate whether your uncertainty makes sense
- ▶ (somewhat) understand how huge deep vision systems work