

PROBABILISTIC INFERENCE AND LEARNING

LECTURE 05

MORE ON GAUSSIAN REGRESSION

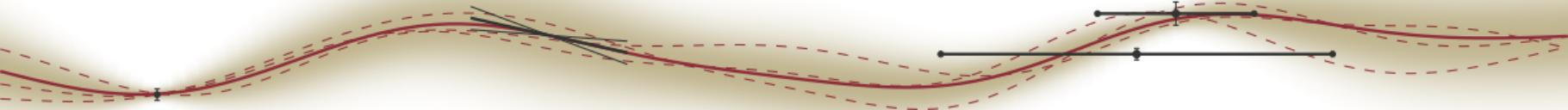
Philipp Hennig

31 October 2018

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

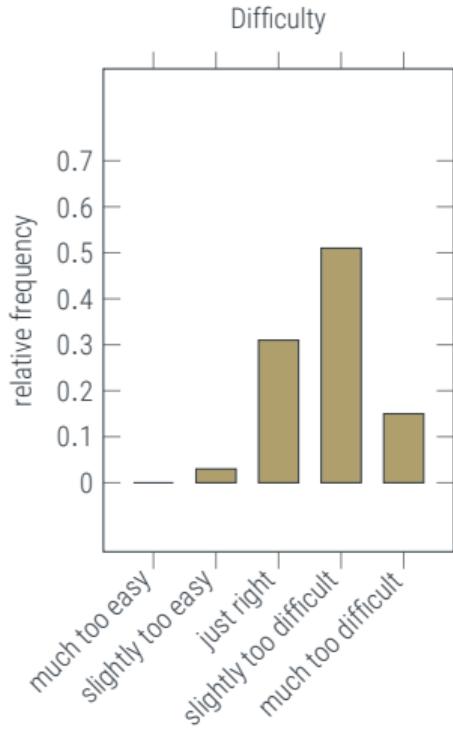
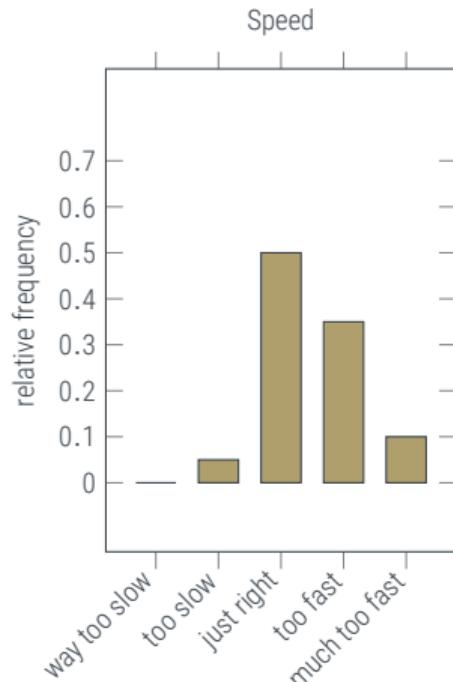
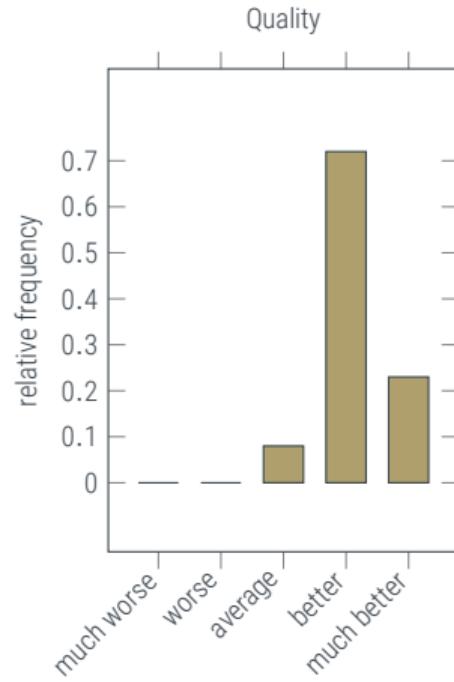


FACULTY OF SCIENCE
DEPARTMENT OF COMPUTER SCIENCE
CHAIR FOR THE METHODS OF MACHINE LEARNING



Last Lecture: Debrief

Feedback dashboard





Last Lecture: Debrief

Detailed Feedback

Things you did not like:

- ♦ long terms on the blackboard
- ♦ I have no python experience
- ♦ "I feel a lot more stupid than every physicist now"
- ♦ black, red and gold are hard to distinguish on the slides

Things you did not understand:

- ♦ weight vs. function space
- ♦ when do we need matrix inverses instead of Cholesky
- ♦ sampling
- ♦ could you please give a preview of the coming lecture at the end?
- ♦ why adding 10^{-6} to the diagonal?

Things you enjoyed:

- ♦ plots / animations (!!)
- ♦ code (!!)
- ♦ pointing out the uncertainty comes for free
- ♦ sampling from Gaussians

Which part of today's lecture did you enjoy **most**?

the sausage of uncertainty





Overview of Lectures so far:

0. Introduction to Reasoning under Uncertainty
 - Probabilities are the mathematical formalization of uncertainty
1. Probabilistic Reasoning
 - Probabilities extend deductive to plausible reasoning. Conditional independence affects complexity
2. Probabilities over Continuous Variables
 - Probability **densities** distribute probability over continuous domains
3. Gaussian Probability Distributions
 - Gaussians map probabilistic inference to **linear algebra**
4. Gaussian Parametric Regression
 - Gaussian's can be used to infer **functions**, using linear models

Today:

- Can we *learn* the features?
- How do we do this in practice?
- Connections to deep learning



Reminder: General Linear Regression

An unbounded abundance of choices for features



Reminder: General Linear Regression

An unbounded abundance of choices for features



Reminder: General Linear Regression

An unbounded abundance of choices for features



Reminder: General Linear Regression

An unbounded abundance of choices for features



Reminder: General Linear Regression

An unbounded abundance of choices for features



Reminder: General Linear Regression

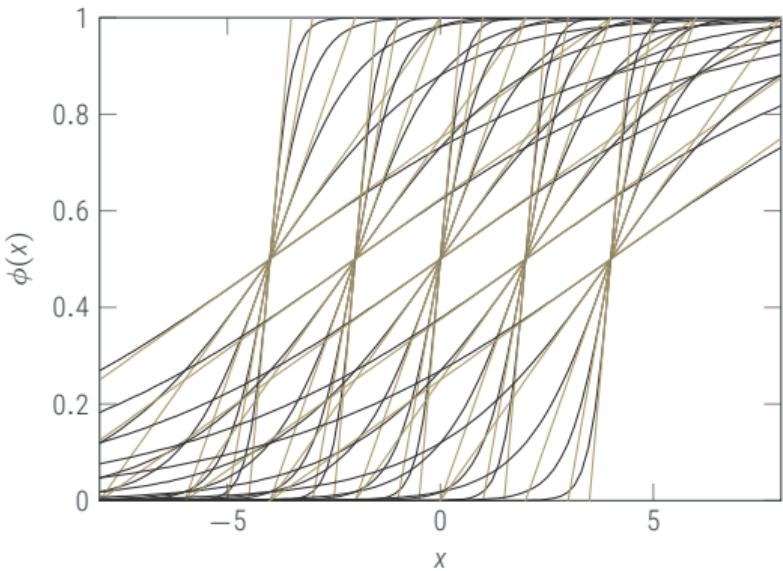
An unbounded abundance of choices for features



Can we Learn the Features?

Hierarchical Bayesian Inference

$$p(w | y, \phi) = \frac{p(y | w, \phi)p(w | \phi)}{p(y | \phi)}$$



- There is an infinite-dimensional space of feature functions to choose from
- Maybe we can restrict to a finite-dimensional sub-space and **search** in there? Say

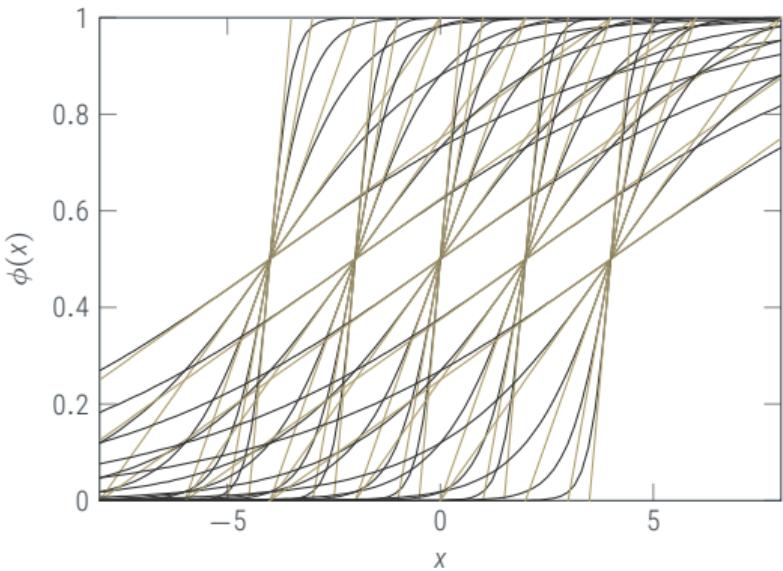
$$\phi_i(x; \boldsymbol{\theta}) = \frac{1}{1 + \exp(-\frac{x-\theta_1}{\theta_2})}$$



Can we Learn the Features?

Hierarchical Bayesian Inference

$$p(w | y, \phi) = \frac{p(y | w, \phi)p(w | \phi)}{p(y | \phi)}$$



- ❖ There is an infinite-dimensional space of feature functions to choose from
- ❖ Maybe we can restrict to a finite-dimensional sub-space and **search** in there? Say

$$\phi_i(x; \boldsymbol{\theta}) = \frac{1}{1 + \exp(-\frac{x-\theta_1}{\theta_2})}$$

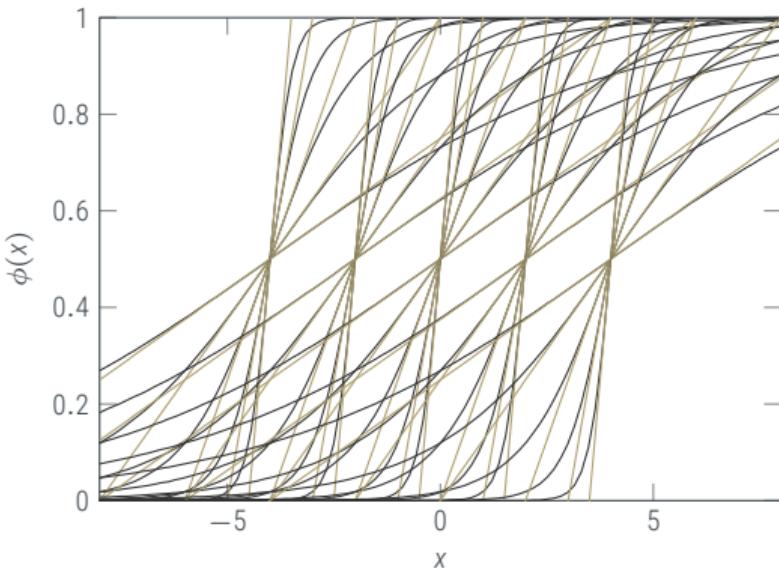
- ❖ θ_1, θ_2 are just unknown parameters!
- ❖ So can we infer them just like w ?



Can we Learn the Features?

Hierarchical Bayesian Inference

$$p(w | y, \phi) = \frac{p(y | w, \phi)p(w | \phi)}{p(y | \phi)}$$



- There is an infinite-dimensional space of feature functions to choose from
- Maybe we can restrict to a finite-dimensional sub-space and **search** in there? Say

$$\phi_i(x; \boldsymbol{\theta}) = \frac{1}{1 + \exp(-\frac{x-\theta_1}{\theta_2})}$$

- θ_1, θ_2 are just unknown parameters!
- So can we infer them just like w ?
- Yes, but not as easily: the likelihood

$$p(y | w, \theta) = \mathcal{N}(y; \phi(x; \boldsymbol{\theta})^T w, \sigma^2)$$

contains a **non-linear** map of $\boldsymbol{\theta}$.

Hierarchical Bayesian Inference

Bayesian model adaptation

$$p(f | y, x, \boldsymbol{\theta}) = \frac{p(y | f, x, \boldsymbol{\theta})p(f |, \boldsymbol{\theta})}{\int p(y | f, x, \boldsymbol{\theta})p(f |, \boldsymbol{\theta}) df} = \frac{p(y | f, x, \boldsymbol{\theta})p(f |, \boldsymbol{\theta})}{p(y | x, \boldsymbol{\theta})}$$

- Model parameters like $\boldsymbol{\theta}$ are also known as hyper-parameters.
- This is largely a computational, practical distinction:

data	are observed	→ condition
variables	are the things we care about	→ full probabilistic treatment
parameters	are the things we have to deal with to get the model right	→ integrate out
hyper-parameters	are the top-level, too expensive to properly infer	→ fit

The **model evidence** in Bayes' Theorem is the (marginal) **likelihood** for the model.

Hierarchical Bayesian Inference

Bayesian model adaptation

$$p(f | y, x, \boldsymbol{\theta}) = \frac{p(y | f, x, \boldsymbol{\theta}) p(f |, \boldsymbol{\theta})}{\int p(y | f, x, \boldsymbol{\theta}) p(f |, \boldsymbol{\theta}) df} = \frac{p(y | f, x, \boldsymbol{\theta}) p(f |, \boldsymbol{\theta})}{p(y | x, \boldsymbol{\theta})}$$

- For Gaussians, die evidence has **analytic form**:

$$\underbrace{\mathcal{N}(y; \phi_X^{\boldsymbol{\theta}^\top} w + b, \Lambda)}_{p(y|f,x,\boldsymbol{\theta})} \cdot \underbrace{\mathcal{N}(w, \mu, \Sigma)}_{p(f)} = \underbrace{\mathcal{N}(w; m_{\text{post}}^{\boldsymbol{\theta}}, V_{\text{post}}^{\boldsymbol{\theta}})}_{p(f|y,x,\boldsymbol{\theta})} \cdot \underbrace{\mathcal{N}(y; \phi_X^{\boldsymbol{\theta}^\top} \mu + b, \phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda)}_{p(y|\boldsymbol{\theta},x)}$$

- BUT:** It's not a linear function of $\boldsymbol{\theta}$, so analytic Gaussian inference is not available!

Computational complexity is *the* principal challenge of probabilistic reasoning.



The Toolbox

Five principal methods for dealing with computational complexity in probabilistic inference

1. Maximum Likelihood (ML) / Maximum A-Posteriori (MAP) estimation:
To estimate θ in $p(D | \theta)$ or $p(\theta | D)$, set $\hat{\theta} = \arg \max_{\theta} p$.
2. ????
3. ????
4. ????
5. ????

Disclaimer: The listed items are neither mutually exclusive nor collectively exhaustive. Some of the methods are intricately interrelated.

ML / MAP in Practice

Finding the "best fit" θ in Gaussian models

[e.g. DJC MacKay, *The evidence framework applied to classification networks*, 1992]

$$\begin{aligned}
 \hat{\boldsymbol{\theta}} &= \arg \max_{\boldsymbol{\theta}} \mathcal{N}(\mathbf{y}; \phi_X^{\boldsymbol{\theta}^\top} \boldsymbol{\mu} + \mathbf{b}, \phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda) \\
 &= \arg \max_{\boldsymbol{\theta}} \log \mathcal{N}(\mathbf{y}; \phi_X^{\boldsymbol{\theta}^\top} \boldsymbol{\mu} + \mathbf{b}, \phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda) \\
 &= \arg \min_{\boldsymbol{\theta}} -\log \mathcal{N}(\mathbf{y}; \phi_X^{\boldsymbol{\theta}^\top} \boldsymbol{\mu} + \mathbf{b}, \phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda) \\
 &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \left(\underbrace{(\mathbf{y} - \phi_X^{\boldsymbol{\theta}^\top} \boldsymbol{\mu})^\top \left(\phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda \right)^{-1} (\mathbf{y} - \phi_X^{\boldsymbol{\theta}^\top} \boldsymbol{\mu})}_{\text{square error}} + \underbrace{\log |\phi_X^{\boldsymbol{\theta}^\top} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda|}_{\text{model complexity / Occam factor}} \right) + \frac{N}{2} \log 2\pi
 \end{aligned}$$



$$\log \left| \phi_X^{\theta^T} \Sigma \phi_X^{\theta} + \Lambda \right|$$

Numquam ponenda est pluralitas sine necessitate.
Plurality must never be posited without necessity.

William of Occam
(1285 (Occam, Surrey)–1349 (Munich, Bavaria))
stained-glass window by Lawrence Lee



What is Model Complexity?

The Occam factor is not always straightforward

$$\log \left| \phi_X^{\boldsymbol{\theta}^T} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda \right|$$

measures model complexity as the “volume” of hypotheses covered by the joint Gaussian distribution.



What is Model Complexity?

The Occam factor is not always straightforward

$$\log \left| \phi_X^{\boldsymbol{\theta}^T} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda \right|$$

measures model complexity as the “volume” of hypotheses covered by the joint Gaussian distribution.



What is Model Complexity?

The Occam factor is not always straightforward

$$\log \left| \phi_X^{\boldsymbol{\theta}^T} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda \right|$$

measures model complexity as the “volume” of hypotheses covered by the joint Gaussian distribution.



What is Model Complexity?

The Occam factor is not always straightforward

$$\log \left| \phi_X^{\boldsymbol{\theta}^T} \Sigma \phi_X^{\boldsymbol{\theta}} + \Lambda \right|$$

measures model complexity as the “volume” of hypotheses covered by the joint Gaussian distribution.



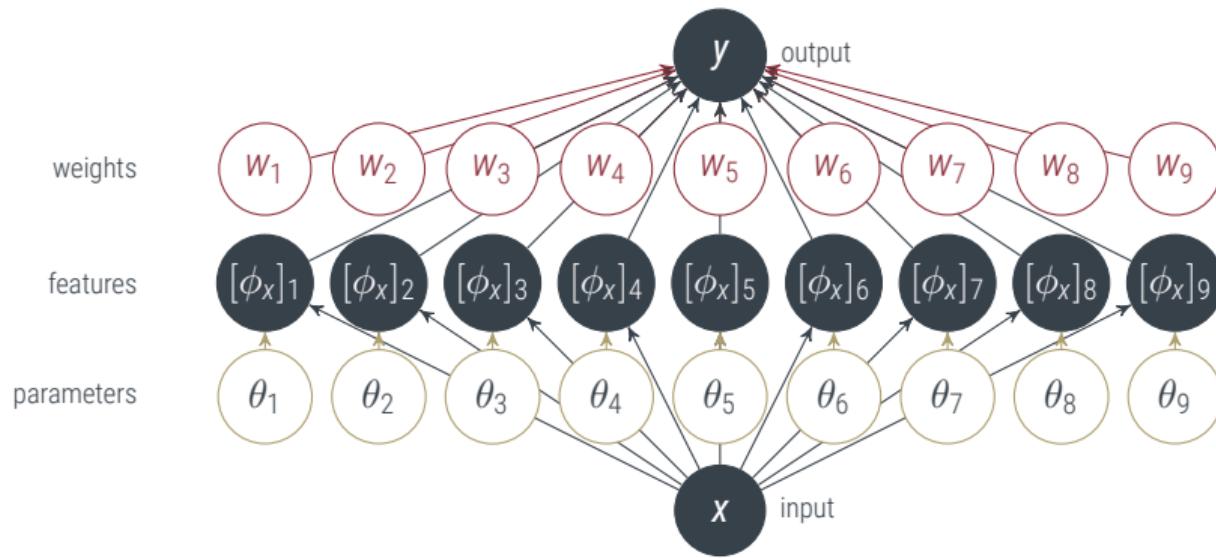
Type II Inference

Fitting a probabilistic model by maximum marginal likelihood



A Structural Observation

Graphical Model



A linear Gaussian regressor is a **single hidden layer** neural network, with quadratic output loss, and fixed input layer. Hyperparameter-fitting corresponds to training the input layer. The usual way to train such network, however, does not include the Occam factor.



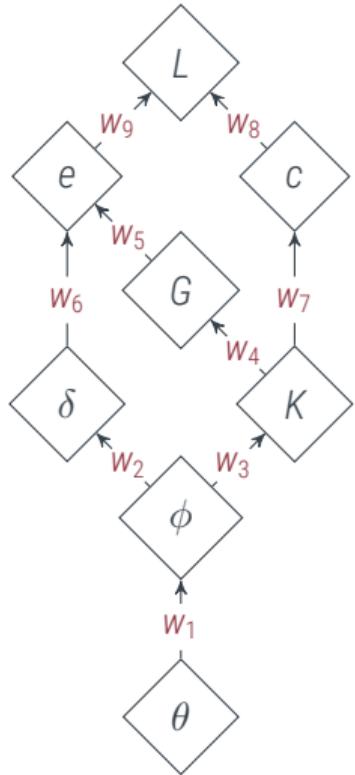
What does the Optimizer need from us?

A bit of algorithmic wizardry

$$L(\boldsymbol{\theta}) = \frac{1}{2} \left(\underbrace{(y - \phi_x^{\boldsymbol{\theta}^\top} \mu)^\top}_{=:e} \underbrace{\left(\underbrace{\phi_x^{\boldsymbol{\theta}^\top} \Sigma \phi_x^{\boldsymbol{\theta}} + \Lambda}_{=:K} \right)^{-1} \underbrace{(y - \phi_x^{\boldsymbol{\theta}^\top} \mu)}_{=: \Delta} + \log \underbrace{\left| \phi_x^{\boldsymbol{\theta}^\top} \Sigma \phi_x^{\boldsymbol{\theta}} + \Lambda \right|}_{=:c}} \right)$$

What does the Optimizer need from us?

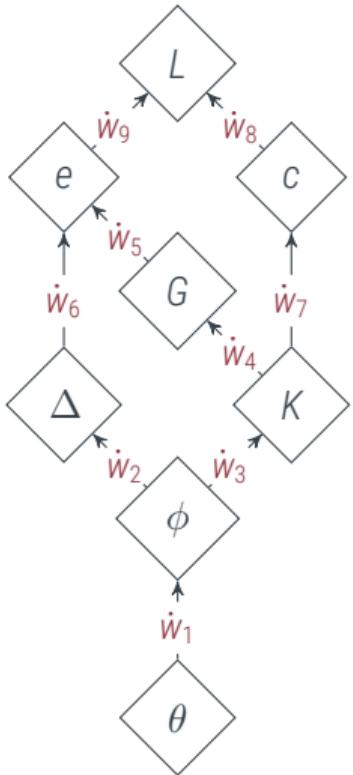
Automatic Differentiation



$$L(\boldsymbol{\theta}) = \frac{1}{2} \left((y - \phi_{\boldsymbol{\theta}}^T \mu)^T \underbrace{\left(\underbrace{\phi_{\boldsymbol{\theta}}^T \Sigma \phi_{\boldsymbol{\theta}}}_{=:K} + \Lambda \right)^{-1}}_{=:G} \underbrace{(y - \phi_{\boldsymbol{\theta}}^T \mu)}_{=:e} + \underbrace{\log |\phi_{\boldsymbol{\theta}}^T \Sigma \phi_{\boldsymbol{\theta}} + \Lambda|}_{=:c} \right)$$

What does the Optimizer need from us?

Automatic Differentiation – Forward Mode

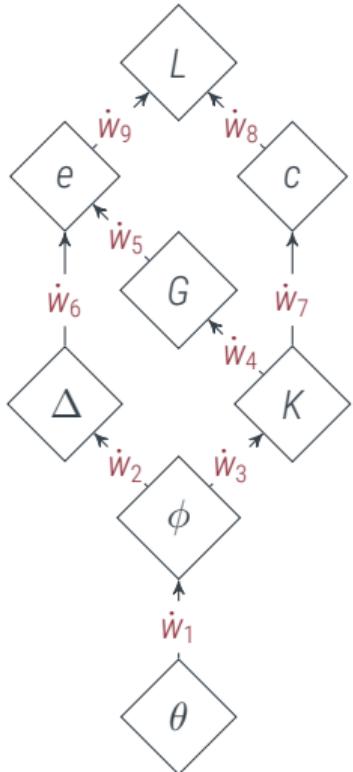


$$L(\boldsymbol{\theta}) = \frac{1}{2} \left((y - \phi_X^\top \mu)^\top \underbrace{\left(\underbrace{\phi_X^\top \Sigma \phi_X}_{=:K} + \Lambda \right)}_{=:G}^{-1} \underbrace{(y - \phi_X^\top \mu)}_{=: \Delta} + \underbrace{\log |\phi_X^\top \Sigma \phi_X + \Lambda|}_{=:c} \right)$$

$$\begin{aligned} \frac{\partial L}{\partial \theta} &= \frac{\partial L}{\partial e} \frac{\partial e}{\partial \theta} + \frac{\partial L}{\partial c} \frac{\partial c}{\partial \theta} = \dot{w}_9 \frac{\partial e}{\partial \theta} + \dot{w}_8 \frac{\partial c}{\partial \theta} = \dot{w}_9 \left(\frac{\partial e}{\partial \Delta} \frac{\partial \Delta}{\partial \theta} + \frac{\partial e}{\partial G} \frac{\partial G}{\partial \theta} \right) + \dot{w}_8 \frac{\partial c}{\partial K} \frac{\partial K}{\partial \theta} \\ &= \dot{w}_9 \left(\dot{w}_6 \frac{\partial \Delta}{\partial \theta} + \dot{w}_5 \frac{\partial G}{\partial \theta} \right) + \dot{w}_8 \dot{w}_7 \frac{\partial K}{\partial \theta} = \dot{w}_9 \left(\dot{w}_6 \frac{\partial \Delta}{\partial \phi} \frac{\partial \phi}{\partial \theta} + \dot{w}_5 \frac{\partial G}{\partial K} \frac{\partial K}{\partial \theta} \right) + \dot{w}_8 \dot{w}_7 \frac{\partial K}{\partial \theta} \\ &= \dot{w}_9 \dot{w}_6 \dot{w}_2 \frac{\partial \phi}{\partial \theta} + (\dot{w}_9 \dot{w}_5 \dot{w}_4 + \dot{w}_8 \dot{w}_7) \frac{\partial K}{\partial \theta} = \left(\dot{w}_9 \dot{w}_6 \dot{w}_2 + (\dot{w}_9 \dot{w}_5 \dot{w}_4 + \dot{w}_8 \dot{w}_7) \frac{\partial K}{\partial \phi} \right) \frac{\partial \phi}{\partial \theta} \\ &= (\dot{w}_9 \dot{w}_6 \dot{w}_2 + (\dot{w}_9 \dot{w}_5 \dot{w}_4 + \dot{w}_8 \dot{w}_7) \dot{w}_3) \frac{\partial \phi}{\partial \theta} \frac{\partial \theta}{\partial \theta} \\ &= (\dot{w}_9 \dot{w}_6 \dot{w}_2 + (\dot{w}_9 \dot{w}_5 \dot{w}_4 + \dot{w}_8 \dot{w}_7) \dot{w}_3) \dot{w}_1 \end{aligned}$$

What does the Optimizer need from us?

Automatic Differentiation – Forward Mode



$$L(\boldsymbol{\theta}) = \frac{1}{2} \left((y - \phi_{\boldsymbol{\theta}}^T \mu)^T \underbrace{\left(\underbrace{\phi_{\boldsymbol{\theta}}^T \Sigma \phi_{\boldsymbol{\theta}}}_{=:K} + \Lambda \right)}_{=:G}^{-1} \underbrace{(y - \phi_{\boldsymbol{\theta}}^T \mu)}_{=:e} + \log \underbrace{\left| \phi_{\boldsymbol{\theta}}^T \Sigma \phi_{\boldsymbol{\theta}} + \Lambda \right|}_{=:c} \right)$$

$$\dot{w}_9 = \frac{\partial L}{\partial e} = 1/2 \quad \dot{w}_8 = \frac{\partial L}{\partial c} = 1/2 \quad [\dot{w}_7]_{ij} = \frac{\partial c}{\partial K_{ij}} = K_{ij}^{-1}$$

$$[\dot{w}_6]_i = \frac{\partial e}{\partial \Delta_i} = 2[G\Delta]_i \quad [\dot{w}_5]_{ij} = \frac{\partial e}{\partial G_{ij}} = \Delta_i \Delta_j \quad [\dot{w}_4]_{ij,k\ell} = \frac{\partial G_{ij}}{\partial K_{k\ell}} = -G_{ik} G_{j\ell}$$

$$[\dot{w}_3]_{ij,ab} = \frac{\partial K_{ij}}{\partial \phi_{ab}} = \delta_{ia} [\Sigma \phi]_{bj} + \delta_{jb} [\Sigma \phi]_{kj}$$

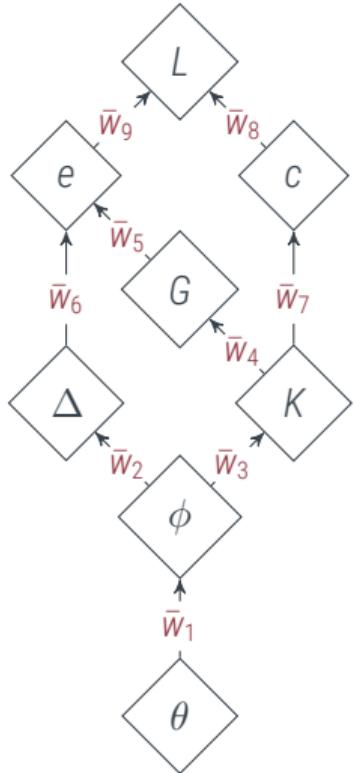
$$[\dot{w}_2]_{i,ab} = \frac{\partial \Delta_i}{\partial \phi_{ab}} = -\delta_{ia} \mu_b \quad [\dot{w}_1]_{ab,\ell} = \frac{\partial \phi_{ab}}{\partial \theta_\ell} = \text{your choice!}$$



[Seppo Linnainmaa, 1970]

What does the Optimizer need from us?

Automatic Differentiation – Backward Mode



$$L(\theta) = \frac{1}{2} \left((y - \phi_x^{\theta \top} \mu)^{\top} \underbrace{\left(\underbrace{\phi_x^{\theta \top} \Sigma \phi_x^{\theta}}_{=:K} + \Lambda \right)}_{=:G}^{-1} \underbrace{(y - \phi_x^{\theta \top} \mu)}_{=: \Delta} + \log \underbrace{\left| \phi_x^{\theta \top} \Sigma \phi_x^{\theta} + \Lambda \right|}_{=:c} \right)$$

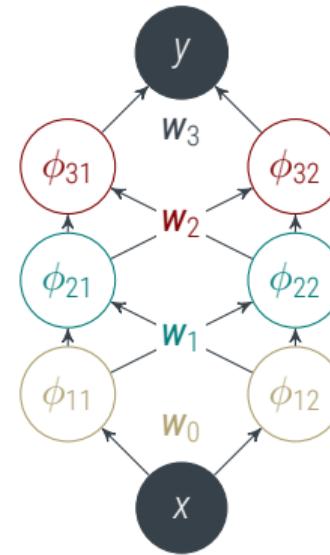
$$\begin{aligned} \frac{\partial L}{\partial \theta} &= \frac{\partial L}{\partial \phi} \frac{\partial \phi}{\partial \theta} =: \bar{w}_1 = \left(\frac{\partial L}{\partial \Delta} \frac{\partial \Delta}{\partial \phi} + \frac{\partial L}{\partial K} \frac{\partial K}{\partial \phi} \right) \frac{\partial \phi}{\partial \theta} =: (\bar{w}_2 + \bar{w}_3) \frac{\partial \phi}{\partial \theta} \\ \bar{w}_2 &= \frac{\partial L}{\partial e} \frac{\partial e}{\partial \Delta} \frac{\partial \Delta}{\partial \phi} =: \bar{w}_6 \frac{\partial \Delta}{\partial \phi} \quad \bar{w}_3 = \left(\frac{\partial L}{\partial G} \frac{\partial G}{\partial K} + \frac{\partial L}{\partial c} \frac{\partial c}{\partial K} \right) \frac{\partial K}{\partial \phi} =: (\bar{w}_4 + \bar{w}_7) \frac{\partial K}{\partial \phi} \\ \bar{w}_4 &= \frac{\partial L}{\partial e} \frac{\partial e}{\partial G} \frac{\partial G}{\partial K} =: \bar{w}_5 \frac{\partial G}{\partial K} \quad \bar{w}_5 = \frac{\partial L}{\partial e} \frac{\partial e}{\partial G} =: \bar{w}_9 \frac{\partial e}{\partial G} \quad \bar{w}_6 = \frac{\partial L}{\partial e} \frac{\partial e}{\partial \Delta} =: \bar{w}_9 \frac{\partial e}{\partial \Delta} \\ \bar{w}_7 &= \frac{\partial L}{\partial c} \frac{\partial c}{\partial K} =: \bar{w}_8 \frac{\partial c}{\partial K} \quad \bar{w}_8 = \bar{w}_9 = 1/2 \end{aligned}$$

$\bar{w}_i = \frac{\partial L}{\partial \text{subgraph}_i}$ are known as *adjoints*. Traverse graph backward to collect the derivative. This is faster than forward-mode for single-output-many-input functions, but requires storing the above structure (known as a *Wengert list*). (cf. "Backpropagation")



Deep Networks

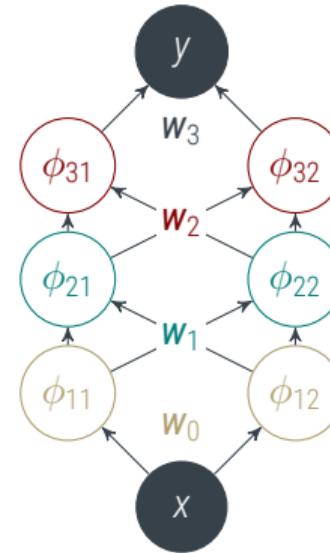
But not Bayesian deep networks



$$\hat{f}(x, W) = \sum_{i=1}^F \phi_{3i}(x, w_{\text{lower}}) w_{3i} = \sum_i \phi_{3i} \left(\sum_j \phi_{2j} \left(\sum_\ell \phi_{1\ell} (w_{0\ell} x) w_{1\ell j} \right) w_{2ji} \right) w_{3i}$$

Deep Networks

But not Bayesian deep networks



$$\hat{f}(x, W) = \arg \min_{W \in \mathbb{R}^D} \|y - \hat{f}(x, W)\|^2 + \alpha^2 \|W\|^2 = \mathcal{L}(W)$$

$$W_{t+1} = W_t + \tau \nabla \mathcal{L}(W)$$



What about Mini-batches?

Learning representations from big data

$$p(w, \theta | y) \propto p(y | w, \phi^\theta) p(w, \theta) = p(w, \theta) \cdot \prod_{i=1}^n p(y_i | w, \phi_i^\theta) = p(w, \theta) \cdot \prod_{i=1}^n \mathcal{N}(y_i; \phi_i^{\theta \top} w, \sigma^2)$$



What about Mini-batches?

Learning representations from big data

$$p(w, \theta | y) \propto p(y | w, \phi^\theta) p(w, \theta) = p(w, \theta) \cdot \prod_{i=1}^n p(y_i | w, \phi_i^\theta) = p(w, \theta) \cdot \prod_{i=1}^n \mathcal{N}(y_i; \phi_i^{\theta \top} w, \sigma^2)$$

$$\arg \max_{w, \theta} p(w, \theta | y) = \arg \min_{w, \theta} -\log p(w, \theta | y)$$

What about Mini-batches?

Learning representations from big data

$$p(w, \theta | y) \propto p(y | w, \phi^\theta) p(w, \theta) = p(w, \theta) \cdot \prod_{i=1}^n p(y_i | w, \phi_i^\theta) = p(w, \theta) \cdot \prod_{i=1}^n \mathcal{N}(y_i; \phi_i^{\theta \top} w, \sigma^2)$$

$$\arg \max_{w, \theta} p(w, \theta | y) = \arg \min_{w, \theta} -\log p(w, \theta | y)$$

$$= \arg \min_{w, \theta} -\log p(w, \theta) - \frac{1}{2\sigma^2} \sum_{i=1}^n \|y_i - \phi_i^{\theta \top} w\|^2 = \arg \min_{w, \theta} r(w, \theta) + \underbrace{\sum_{i=1}^n \ell_2(y_i; \theta, w)}_{\mathcal{L}(\theta, w)}$$



What about Mini-batches?

Learning representations from big data

$$p(w, \theta | y) \propto p(y | w, \phi^\theta) p(w, \theta) = p(w, \theta) \cdot \prod_{i=1}^n p(y_i | w, \phi_i^\theta) = p(w, \theta) \cdot \prod_{i=1}^n \mathcal{N}(y_i; \phi_i^{\theta \top} w, \sigma^2)$$

$$\arg \max_{w, \theta} p(w, \theta | y) = \arg \min_{w, \theta} -\log p(w, \theta | y)$$

$$= \arg \min_{w, \theta} -\log p(w, \theta) - \frac{1}{2\sigma^2} \sum_{i=1}^n \|y_i - \phi_i^{\theta \top} w\|^2 = \arg \min_{w, \theta} r(w, \theta) + \underbrace{\sum_{i=1}^n \ell_2(y_i; \theta, w)}_{\mathcal{L}(\theta, w)}$$

$$\approx \arg \min_{w, \theta} -\log p(w, \theta) - \frac{1}{2\sigma^2} \frac{n}{b} \sum_{\beta=1}^b \|y_\beta - \phi_\beta^{\theta \top} w\|^2$$



What about Mini-batches?

Learning representations from big data

$$p(w, \theta | y) \propto p(y | w, \phi^\theta) p(w, \theta) = p(w, \theta) \cdot \prod_{i=1}^n p(y_i | w, \phi_i^\theta) = p(w, \theta) \cdot \prod_{i=1}^n \mathcal{N}(y_i; \phi_i^{\theta \top} w, \sigma^2)$$

$$\arg \max_{w, \theta} p(w, \theta | y) = \arg \min_{w, \theta} -\log p(w, \theta | y)$$

$$= \arg \min_{w, \theta} -\log p(w, \theta) - \frac{1}{2\sigma^2} \sum_{i=1}^n \|y_i - \phi_i^{\theta \top} w\|^2 = \arg \min_{w, \theta} r(w, \theta) + \underbrace{\sum_{i=1}^n \ell_2(y_i; \theta, w)}_{\mathcal{L}(\theta, w)}$$

$$\approx \arg \min_{w, \theta} -\log p(w, \theta) - \frac{1}{2\sigma^2} \frac{n}{b} \sum_{\beta=1}^b \|y_\beta - \phi_\beta^{\theta \top} w\|^2 \sim \mathcal{N}(r + \mathcal{L}(\theta, w), \mathcal{O}(b^{-1}))$$



Connections and Differences

Bayesian and Deep Learning

- ⊕ MAP inference does not capture **uncertainty** on parameters:
 - ⊕ no posterior uncertainty from not fully identified parameters
 - ⊕ no model capacity control in the evidence term
- ⊕ A linear Gaussian regressor is a **single hidden layer** neural network, with quadratic output loss, and fixed input layer (deep networks can of course be treated in the same way).
Hyperparameter-fitting corresponds to training the input layer. The usual way to train such network, however, does not include the Occam factor. Data sub-sampling can be used just as in other areas to speed up computations at the cost of reduced computational precision.
- ⊕ All worries one may have about fitting or hand-picking features for Bayesian regression also apply to deep learning. By highlighting assumptions and priors, the probabilistic view forces us to address many problems directly, rather than obscuring them with notation and intuitions.
- ⊕ **Automatic Differentiation (AD)** is a algorithmic tool that is just as helpful for Bayesian inference as it is for deep learning

It is possible to construct a point estimate for a Bayesian model, and to construct full posteriors for deep networks. The two domains are not separate, they are just different mental scaffolds. If you're hoping for a theory of deep learning, probability theory is a primary contender.



Summary:

- The features used for Gaussian linear regression can be **learnt** by **hierarchical Bayesian Inference**
- This is usually **intractable**. Instead, **approximate inference** methods are used
- For example, **maximal a-posteriori probability (MAP)** inference fits a point-estimate for feature parameters
- MAP inference is an **optimization** problem, and can thus be performed in the same way as other optimization-based ML approaches, including deep learning. That is, using the same optimizers (e.g. stochastic gradient descent), the same automatic differentiation frameworks (e.g. TensorFlow / pyTorch, etc.) and the same data subsampling techniques.

The different viewpoints (probabilistic / statistical / empirical ("deep")) on Machine Learning often overlap and inform each other. Understanding of Bayesian linear (Gaussian) regression can help us build a better intuition for deep learning, too.

Next lecture: Instead of *learning* a few features, sometimes we can get away with using *infinitely many* features.