

# 02477 – Bayesian Machine Learning: Lecture 3

Michael Riis Andersen

Technical University of Denmark,  
DTU Compute, Department of Applied Math and Computer Science

# Outline

- 1 Recap: Probabilistic machine learning
- 2 Recap: Multivariate Gaussian distributions
- 3 Recap: Linear regression and supervised learning
- 4 Bayesian linear regression
- 5 The posterior predictive distribution
- 6 Dealing with hyperparameters

## Recap: Probabilistic machine learning

# A probabilistic perspective on making predictions

## Product rule

$$p(\mathbf{a}, \mathbf{b}) = p(\mathbf{b}|\mathbf{a})p(\mathbf{a})$$

## Sum rule

$$p(\mathbf{b}) = \int p(\mathbf{a}, \mathbf{b})d\mathbf{a}$$

## Conditional

$$p(\mathbf{a}|\mathbf{b}) = \frac{p(\mathbf{a}, \mathbf{b})}{p(\mathbf{b})}$$

## Conditional independence

$$p(\mathbf{a}, \mathbf{b}|\mathbf{c}) = p(\mathbf{a}|\mathbf{c})p(\mathbf{b}|\mathbf{c})$$

**Goal:** Given some data  $y$ , what can we say about a new observation  $y^*$ ?

- Step 1: Formulate *joint distribution* for *all variables* of interests

$$p(y^*, y, \theta) = p(y^*, y|\theta)p(\theta) = p(y^*|\theta)p(y|\theta)p(\theta)$$

- Step 2: *Condition* on the *observed data*  $y$

$$p(y^*, \theta|y) = \frac{p(y^*, y, \theta)}{p(y)} = \frac{p(y^*|\theta)p(y|\theta)p(\theta)}{p(y)}$$

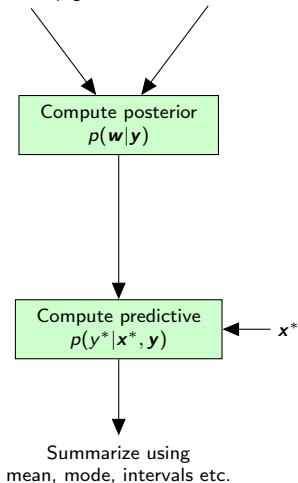
- Step 3: *Marginalize* out parameter  $\theta$  using the *sum rule* to get the *posterior predictive distribution*

$$p(y^*|y) = \int p(y^*, \theta|y)d\theta = \int \frac{p(y^*|\theta)p(y|\theta)p(\theta)}{p(y)}d\theta = \int p(y^*|\theta)p(\theta|y)d\theta = \mathbb{E}_{p(\theta|y)} [p(y^*|\theta)]$$

- **Key take-away:** To reason about  $y^*$  *given*  $y$ , we need to *average the likelihood* for  $y^*$  wrt. to the *posterior distribution*  $p(\theta|y)$ .

## Bayesian inference for supervised learning

$$\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N \quad p(\mathbf{w}, \mathbf{y}) = p(\mathbf{y}|\mathbf{w})p(\mathbf{w})$$



- Same principles for linear regression, logistic regression, neural networks etc. etc.

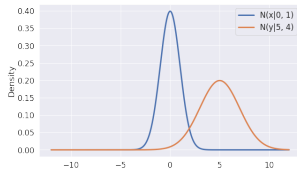
## Recap: Multivariate Gaussian distributions

# Univariate normal distribution

- The *normal distribution* (also known as the Gaussian)

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

- Two parameters:  $\mu = \mathbb{E}[x]$  and  $\sigma^2 = \mathbb{V}[x]$
- Widely due to Central limit theorems, maximum entropy principle, relation to least squares minimization, nice mathematical properties



- Closed under affine transformations

$$x \sim \mathcal{N}(m, v) \quad \Rightarrow \quad a + bx \sim \mathcal{N}(a + bm, b^2v)$$

$$x \sim \mathcal{N}(0, 1)$$

$$y = 5 + 2x$$

- Let  $x \sim \mathcal{N}(m_x, v_x)$  and  $y \sim \mathcal{N}(m_y, v_y)$  for  $x \perp y$ , then

$$x + y \sim \mathcal{N}(m_x + m_y, v_x + v_y)$$

- Functional form: The *logarithm of a Gaussian density* is a *second order polynomial*

$$\ln \mathcal{N}(x|\mu, \sigma^2) = -\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x + K$$

## The functional form of a Gaussian distribution I

- Recall we discussed the *functional form* of a Beta distribution
- Let's derive the functional form of the Gaussian

$$\ln \mathcal{N}(x|\mu, \sigma^2) = \ln \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(x - \mu)^2}{2\sigma^2} \right) \right]$$



## The functional form of a Gaussian distribution II: example

- Take-away: the *functional form* of a univariate Gaussian is

$$\ln \mathcal{N}(x|\mu, \sigma^2) = -\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x + K$$

- Example: suppose we are given the following log density for a random variable  $x$

$$\ln p(x) = -\frac{1}{4}x^2 + 2x + K$$

- We recognize the 2nd order polynomial and conclude that  $p(x)$  must be Gaussian
- We determine the *variance* by matching the coefficient for 2nd order term

$$-\frac{1}{2\sigma^2} = -\frac{1}{4} \quad \Rightarrow \quad \frac{1}{\sigma^2} = \frac{1}{2} \quad \Rightarrow \quad \sigma^2 = 2$$

- We determine the *mean* by matching coefficient for 1st order term

$$\frac{\mu}{\sigma^2} = 2 \quad \Rightarrow \quad \mu = 4$$

- Therefore, we conclude  $p(x) = \mathcal{N}(x|4, 2)$ .

# The multivariate normal distribution

- The *multivariate normal distribution*

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

- Two parameters:  $\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}]$  and  $\boldsymbol{\Sigma} = \text{cov}[\mathbf{x}]$

- Covariance matrix for  $D = 2$

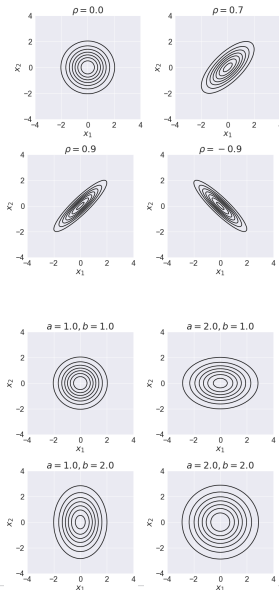
$$\boldsymbol{\Sigma} = \begin{bmatrix} a & c \\ c & b \end{bmatrix}$$

- Correlation coefficient  $\rho = \frac{c}{\sqrt{ab}}$

- Let  $\mathbf{x} \sim \mathcal{N}(\mathbf{m}_x, \mathbf{V}_x)$  and  $\mathbf{y} \sim \mathcal{N}(\mathbf{m}_y, \mathbf{V}_y)$  for  $\mathbf{x} \perp \mathbf{y}$ , then

$$\mathbf{a} + \mathbf{B}\mathbf{x} \sim \mathcal{N}(\mathbf{a} + \mathbf{B}\mathbf{m}_x, \mathbf{B}\mathbf{V}_x\mathbf{B}^T)$$

$$\mathbf{x} + \mathbf{y} \sim \mathcal{N}(\mathbf{m}_x + \mathbf{m}_y, \mathbf{V}_x + \mathbf{V}_y)$$



## The functional form of multivariate Gaussians

- Consider now the log density, focussing only on terms dependent on  $\mathbf{x}$ .

$$\begin{aligned}\ln p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \ln \left[ (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left( -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right) \right] \\&= -\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \\&= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) + \text{constant} \\&= -\frac{1}{2}(\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) + \text{constant} \\&= -\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \text{constant}\end{aligned}$$

- Key take-aways

1. Every time we encounter a distribution with a *quadratic log density*, it must be a Gaussian distribution (if  $\boldsymbol{\Sigma}$  is a valid covariance matrix)
2. We can *match coefficients* of first and second order term to *determine mean and covariance*

Recap: Linear regression and supervised learning

# Supervised learning: linear regression

- Dataset  $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$ 
  - Input features:  $\mathbf{x}_i \in \mathbb{R}^D$
  - Targets:  $y_i \in \mathbb{R}$

- Additive noise models

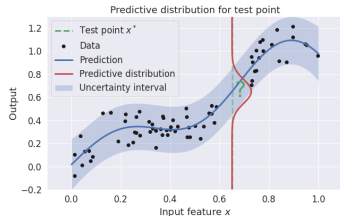
$$y_i = f(\mathbf{x}_i | \mathbf{w}) + \epsilon_i$$

- Linear models are *linear wrt. parameters*, not data!

$$f(\mathbf{x} | \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_M x_M = \mathbf{w}^T \mathbf{x}$$

- *Non-linear* feature extractors  $\phi(\cdot)$  (basis functions)

$$f(\mathbf{x} | \mathbf{w}) = \sum_{j=0}^M w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$



## Linear regression: the probabilistic model

- The *predictive distribution* of  $y^*$  given  $\mathbf{x}^*$  and the data  $\mathcal{D}$  is our goal

$$p(y^*|\mathcal{D}, \mathbf{x}^*)$$

- Model for the "signal"

$$f(\mathbf{x}_i|\mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}_i)$$

- The gaussian noise  $\epsilon_i$  is assumed to be *independent and identically distributed (i.i.d)*

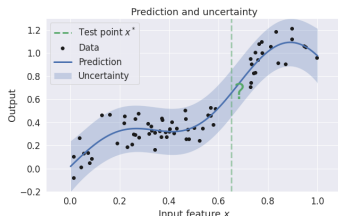
$$y_i = f(\mathbf{x}_i|\mathbf{w}) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

- The *likelihood* for the  $i$ 'th data point

$$p(y_i|\mathbf{x}_i, \mathbf{w}, \sigma^2) = \mathcal{N}(y_i|\mathbf{w}^T \phi(\mathbf{x}_i), \sigma^2)$$

- Using the maximum likelihood solution as a *plug-in* estimator

$$p(y^*|\mathcal{D}, \mathbf{x}^*) = \mathcal{N}(y_i|\hat{\mathbf{w}}_{\text{MLE}}^T \phi(\mathbf{x}_*), \sigma^2)$$



## Estimating the parameters using maximum likelihood

- Given a dataset  $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$ , the *likelihood* for dataset is

$$p(\mathbf{y}|\mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(y_n | f(\mathbf{x}_n | \mathbf{w}), \sigma^2)$$

- Taking the logarithm and using  $f(\mathbf{x}_n, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}_n)$

$$\begin{aligned} \ln p(\mathbf{y}|\mathbf{w}, \sigma^2) &= \sum_{n=1}^N \ln \mathcal{N}(y_n | \mathbf{w}^T \phi(\mathbf{x}_n), \sigma^2) \\ &= \sum_{n=1}^N \left[ -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 \right] \\ &= -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 \end{aligned}$$

- Maximum likelihood estimator  $\hat{\mathbf{w}}_{\text{MLE}}$  is equivalent to minimizing sum-of-squares error

$$\hat{\mathbf{w}}_{\text{MLE}} = \left( \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{y} \quad (\text{Normal equations})$$

$$\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{\mathbf{w}}_{\text{MLE}}^T \phi(\mathbf{x}_n))^2$$

## Example: Polynomial regression using maximum likelihood I

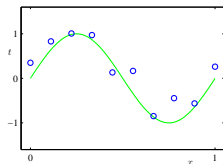
Example from Bishop

### ■ Polynomial basis functions

$$f(x|\mathbf{w}) = w_0 + w_1x + w_2x^2 + \cdots + w_Mx^M = \mathbf{w}^T \boldsymbol{\phi}(x)$$

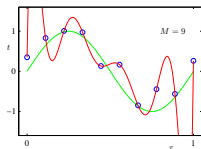
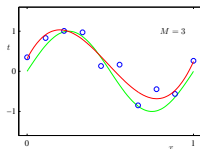
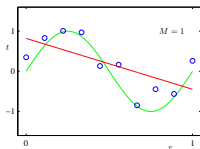
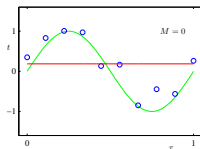
### ■ Feature transformations

$$\boldsymbol{\phi}(x) = [1 \quad x \quad x^2 \quad \cdots \quad x^M]^T$$



### ■ $M$ controls the *model complexity*: Underfitting vs overfitting

### ■ *Model selection*: How to choose $M$ ?





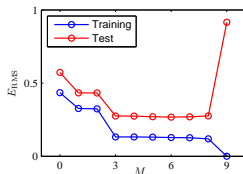
## Example: Polynomial regression using maximum likelihood II

### ■ Cross-validation

Split data into training and test sets

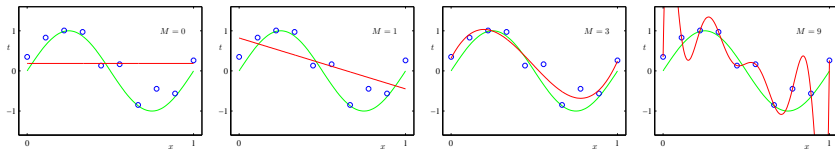
### ■ Overfitting (low training error, high test error)

As the function become more flexible we start to fit the noise in the data



### ■ "Underfitting" (high training error, high test error)

When the function is not sufficiently flexible to fit the data



## Example: Polynomial regression using maximum likelihood III

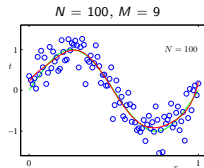
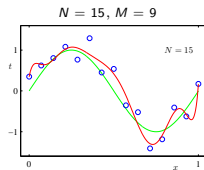
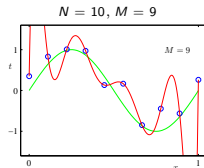
- *The optimal model complexity depends on the amount of data*

The more data the more flexible model we can "afford" to fit

- *Regularization: Controlling the effective model complexity*

We can use regularization to control the effect model complexity when we have limited data

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
$w_0^*$	0.19	0.82	0.31	0.35
$w_1^*$		-1.27	7.99	232.37
$w_2^*$			-25.43	-5321.83
$w_3^*$			17.37	48568.31
$w_4^*$				-231639.30
$w_5^*$				640042.26
$w_6^*$				-1061800.52
$w_7^*$				1042400.18
$w_8^*$				-557682.99
$w_9^*$				125201.43



## Regularized least squares

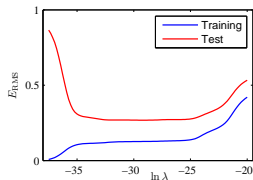
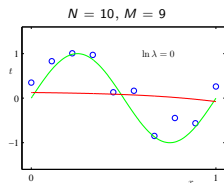
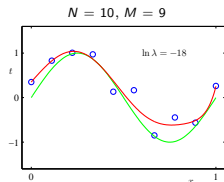
- Recall: *Maximum likelihood* is equivalent *minimizing sum-of-squares error*

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y_n - f(x_n|\mathbf{w}))^2$$

- Adding *penalty term* to prevent weights from becoming too large

$$\tilde{E}_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y_n - f(x_n|\mathbf{w}))^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- What happens when  $\lambda = 0$ ?  $\lambda \rightarrow \infty$ ?
- Many names: Ridge regression, shrinkage, weight decay
- Regularization parameter  $\lambda$  controls the *effective complexity*



## Bayesian linear regression

# Bayesian Linear regression: motivation

## ■ *Overfitting*

Maximum likelihood can be problematic for flexible models

## ■ *Controlling model complexity*

Limiting number of basis functions and/or regularization?

## ■ *Model selection*

How to choose the optimal value of  $\lambda$ ?

## ■ *Cross-validation*

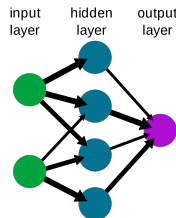
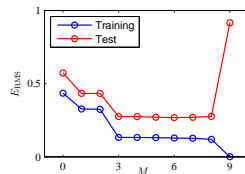
Training + validation/development + test

## ■ *Bayesian methods*

- less prone to overfitting
- can (often) adapt model complexity automatically

## ■ *Applications in modern machine learning*

1. Small datasets
2. Transfer learning
3. Component in more complex models
4. Simple uncertainty quantification for neural networks



# Bayesian Linear regression: prior and likelihood

- Simplified set-up: assuming  $\sigma^2$  is fixed and known, then Bayes' rule states

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

- We already know the *likelihood*

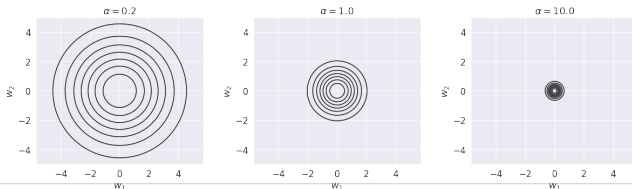
$$p(\mathbf{y}|\mathbf{x}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{w}^T \phi(\mathbf{x}_n), \sigma^2) = \mathcal{N}(\mathbf{y} | \Phi \mathbf{w}, \sigma^2 \mathbf{I})$$

- *The marginal likelihood* is the denominator in Bayes's theorem and is independent of  $\mathbf{w}$

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{w})p(\mathbf{w})d\mathbf{w} = \mathbb{E}_{p(\mathbf{w})} [p(\mathbf{y}|\mathbf{w})]$$

- The multivariate normal distribution is a *conjugate prior* for the  $\mathbf{w}$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$$



## Bayesian Linear regression: The MAP estimator

- The *posterior distribution* of the weights  $\mathbf{w}$  given the data  $\mathbf{y}$  is given by

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{y})} = \frac{\mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2\mathbf{I})\mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})}{p(\mathbf{y})}$$

- Let's look at the *maximum a posteriori* (MAP) estimate:  $\hat{\mathbf{w}}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathbf{y})$

$$p(\mathbf{w}|\mathbf{y}) \propto \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2\mathbf{I})\mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$$

- Taking the logarithm

$$\begin{aligned}\ln p(\mathbf{w}|\mathbf{y}) &\propto \ln \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2\mathbf{I}) + \ln \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) \\ &= -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 - \frac{D}{2} \ln(2\pi\alpha^{-1}) - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \\ &= -\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 - \frac{\alpha}{2} \sum_{i=1}^D w_i^2 + \text{constant}\end{aligned}$$

- The *mode of the posterior (MAP)* is equivalent to ridge regression with  $\lambda = \frac{\alpha}{\beta}$  and to maximum likelihood when  $\alpha \rightarrow 0$

## Deriving the posterior distribution of the weights

Recall the functional form of a generic multivariate Gaussian  $\mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S})$

$$\ln \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S}) = -\frac{1}{2} \mathbf{w}^T \mathbf{S}^{-1} \mathbf{w} + \mathbf{m}^T \mathbf{S}^{-1} \mathbf{w} + \text{constant}$$

- We focus on term that depends on  $\mathbf{w}$ . From the previous slide, we have

$$\begin{aligned} \ln p(\mathbf{w}|\mathbf{y}) &= -\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \phi(x_n))^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{constant} \\ &= -\frac{\beta}{2} (\mathbf{y} - \Phi \mathbf{w})^T (\mathbf{y} - \Phi \mathbf{w}) - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{constant} \\ &= -\frac{\beta}{2} (\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \Phi \mathbf{w} + \mathbf{w} \Phi^T \Phi \mathbf{w}) - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{constant} \\ &= -\frac{\beta}{2} (-2\mathbf{y}^T \Phi \mathbf{w} + \mathbf{w} \Phi^T \Phi \mathbf{w}) - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{constant} \\ &= -\frac{1}{2} \mathbf{w}^T (\beta \Phi^T \Phi + \alpha \mathbf{I}) \mathbf{w} + \beta \mathbf{y}^T \Phi \mathbf{w} + \text{constant} \end{aligned}$$

- Equating coefficients for the second order term

$$\mathbf{S}^{-1} = \beta \Phi^T \Phi + \alpha \mathbf{I} \quad \Longleftrightarrow \quad \mathbf{S} = (\beta \Phi^T \Phi + \alpha \mathbf{I})^{-1}$$

- Equating coefficients for the first order term

$$\mathbf{m}^T \mathbf{S}^{-1} = \beta \mathbf{y}^T \Phi \quad \Longleftrightarrow \quad \mathbf{m} = \beta \mathbf{S} \Phi^T \mathbf{y}$$



## Bayesian linear regression model: the key equations

- Given design matrix  $\Phi \in \mathbb{R}^{N \times D}$  and observations  $\mathbf{y} \in \mathbb{R}^N$ :

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}) \quad (\text{prior})$$

$$p(\mathbf{y} | \mathbf{w}) = \mathcal{N}(\mathbf{y} | \Phi \mathbf{w}, \sigma^2 \mathbf{I}) \quad (\text{likelihood})$$

$$p(\mathbf{w} | \mathbf{y}) = \mathcal{N}(\mathbf{w} | \mathbf{m}, \mathbf{S}) \quad (\text{posterior})$$

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \sigma^2 \mathbf{I} + \alpha^{-1} \Phi \Phi^T) \quad (\text{marginal likelihood})$$

- The *posterior parameters* are given by (using  $\beta \equiv \frac{1}{\sigma^2}$ )

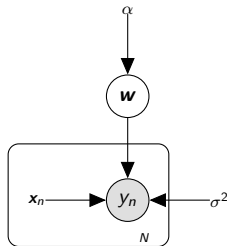
$$\mathbf{m} = \beta \mathbf{S} \Phi^T \mathbf{y}$$

$$\mathbf{S} = \left( \alpha \mathbf{I} + \beta \Phi^T \Phi \right)^{-1}$$

- Two hyperparameters*

$\alpha$ : prior precision of the regression weights

$\beta$ : precision of the measurements



## Linear Gaussian-systems in general (see Section 3.3 in Murphy1)

- For *linear* systems: the Gaussian distribution is *conjugate* to itself
- The *posterior* for a *linear* Gaussian model with Gaussian prior is also *Gaussian*

$$p(\mathbf{y}|\mathbf{z}) = \mathcal{N}(\mathbf{y}|\mathbf{W}\mathbf{z} + \mathbf{b}, \Sigma_y) \qquad p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_z, \Sigma_z)$$

- The *joint* distribution  $p(\mathbf{z}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{z} \\ \mathbf{y} \end{bmatrix} \middle| \boldsymbol{\mu}, \boldsymbol{\Sigma}\right)$

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_z \\ \mathbf{W}\boldsymbol{\mu}_z + \mathbf{b} \end{bmatrix} \qquad \boldsymbol{\Sigma} = \begin{bmatrix} \Sigma_z & \Sigma_z \mathbf{W}^T \\ \mathbf{W}\Sigma_z & \Sigma_y + \mathbf{W}\Sigma_z \mathbf{W}^T \end{bmatrix}$$

- The *posterior* distribution of  $\mathbf{z}$  given  $\mathbf{y}$

$$\begin{aligned} p(\mathbf{y}|\mathbf{z}) &= \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{z|y}, \Sigma_{z|y}) \\ \Sigma_{z|y}^{-1} &= \Sigma_z^{-1} + \mathbf{W}^T \Sigma_y^{-1} \mathbf{W} \\ \boldsymbol{\mu}_{z|y} &= \Sigma_{z|y} \left[ \mathbf{W}^T \Sigma_y^{-1} (\mathbf{y} - \mathbf{b}) + \Sigma_z^{-1} \boldsymbol{\mu}_z \right] \end{aligned}$$

- The *marginal* distribution  $\mathbf{y}$

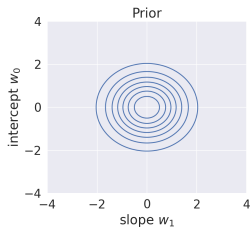
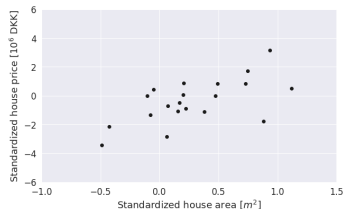
$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{W}\boldsymbol{\mu}_z + \mathbf{b}, \Sigma_y + \mathbf{W}\Sigma_z \mathbf{W}^T)$$

## Example

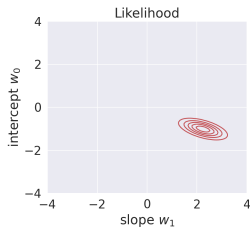
- Simple linear model for fictive house prices

$$f(x|\mathbf{w}) = w_0 + w_1x$$

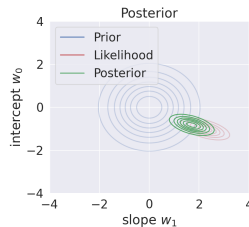
- The posterior summarizes our beliefs about the parameters after seeing the data



$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|0, \alpha^{-1}\mathbf{I})$$



$$p(\mathbf{y}|\mathbf{w}) = \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2\mathbf{I})$$



$$p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S})$$

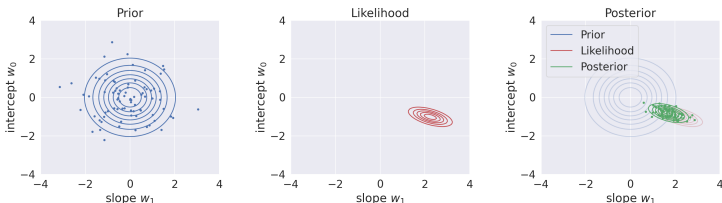
## Posterior inference

$$\text{house price} = w_0 + w_1 \cdot \text{area}$$

- Question: Are larger areas associated with bigger house prices?
- We can calculate various probabilities of interest directly from the posterior (analytically or via sampling), e.g.

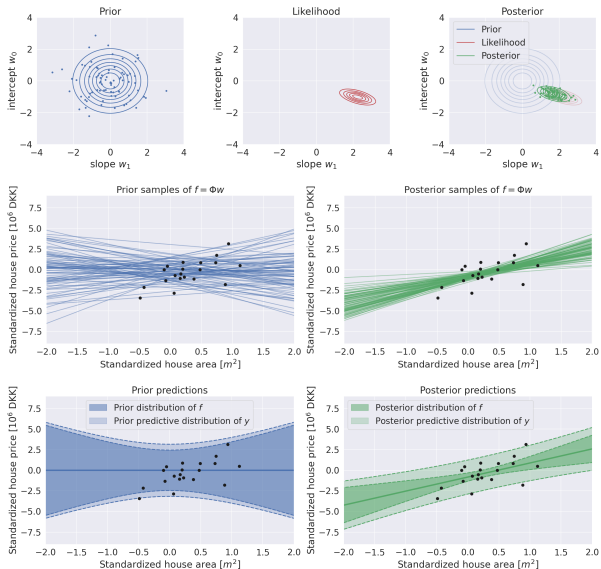
$$p(w_1 > 0 | y) \approx 0.99$$

- No need to remember whether to should use t-tests, F-tests,  $\chi^2$ -tests etc



## The posterior predictive distribution

# Prior and posterior predictive distributions



# Posterior Predictive distributions I

- The *posterior distribution* is  $p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S})$  with

$$\mathbf{m} = \beta \mathbf{S} \Phi^T \mathbf{y}$$

$$\mathbf{S} = (\alpha \mathbf{I} + \beta \Phi \Phi^T)^{-1}$$

- When making predictions  $\mathbf{x}^*$  using Bayesian methods, we *average over all possible parameters values weighted by the posterior*

$$f(\mathbf{x}^*|\mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}^*)$$

$$y(\mathbf{x}^*|\mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}^*) + \epsilon$$

- First, we write the likelihood corresponding to the new input  $\mathbf{x}^*$

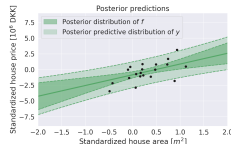
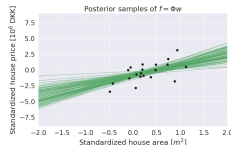
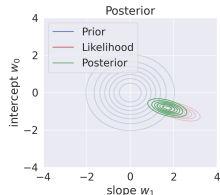
$$p(y^*|\mathbf{x}^*, \mathbf{w}) = \mathcal{N}(y^*|\mathbf{w}^T \phi(\mathbf{x}^*), \sigma^2)$$

- .. and then we integrate with respect to the posterior distribution

$$p(y^*|\mathbf{y}, \mathbf{x}^*) = \int p(y^*, \mathbf{w}|\mathbf{y}, \mathbf{x}^*) d\mathbf{w} \quad (\text{Sum rule})$$

$$= \int p(y^*|\mathbf{x}^*, \mathbf{w}) p(\mathbf{w}|\mathbf{y}) d\mathbf{w} \quad (\text{Product rule})$$

- This is called the *the posterior predictive distribution*



## Posterior Predictive distributions II: 5 minutes exercise

- Our linear model with additive noise is given by (short hand notation:  $\phi = \phi(\mathbf{x})$ )

$$y = \mathbf{w}^T \phi + \epsilon \quad \text{for} \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- Let  $\mathbf{x} \sim \mathcal{N}(\mathbf{m}_x, \mathbf{V}_x)$  and  $\mathbf{y} \sim \mathcal{N}(\mathbf{m}_y, \mathbf{V}_y)$  for  $\mathbf{x} \perp \mathbf{y}$ , then

$$\mathbf{a} + \mathbf{B}\mathbf{x} \sim \mathcal{N}(\mathbf{a} + \mathbf{B}\mathbf{m}_x, \mathbf{B}\mathbf{V}_x\mathbf{B}^T)$$

$$\mathbf{x} + \mathbf{y} \sim \mathcal{N}(\mathbf{m}_x + \mathbf{m}_y, \mathbf{V}_x + \mathbf{V}_y)$$

**Questions:** If our posterior of the weights is given by  $p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S})$ , ...

1. What is the posterior distribution of  $f = \mathbf{w}^T \phi$ ?

Hints:

$$\mathbf{w}^T \phi = 0 + \phi^T \mathbf{w}$$

2. What is the posterior (predictive) distribution of  $y = \mathbf{w}^T \phi + \epsilon$ ?



## Posterior Predictive distributions III

- Generally, we need to solve these integrals analytically or by sampling, but here we can use our calculations from earlier

$$p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S})$$

- The *posterior predictive distribution* of  $y^* = \mathbf{w}^T \phi(\mathbf{x}^*) + \epsilon$  is

$$y(\mathbf{x}^*) \sim \mathcal{N}(m_*, \sigma_*^2)$$

where

$$m_* = \mathbf{m}^T \phi(\mathbf{x}^*)$$

$$\sigma_*^2 = \phi(\mathbf{x}^*)^T \mathbf{S} \phi(\mathbf{x}^*) + \sigma^2$$

- The first term in  $\sigma_*^2$  is due to parameter uncertainty (*epistemic/reducible*) and the second term is due to measurement noise (*aleatoric/irreducible*)
- The term  $\phi(\mathbf{x}_*)^T \mathbf{S} \phi(\mathbf{x}_*)$  is the *posterior uncertainty projected to data space*
- What happens to the predictive variance when  $N \rightarrow \infty$ ?

## Quiz time!

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{y})} \quad (\textit{Bayes' rule})$$

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{w})p(\mathbf{w})d\mathbf{w} \quad (\textit{marginal likelihood})$$

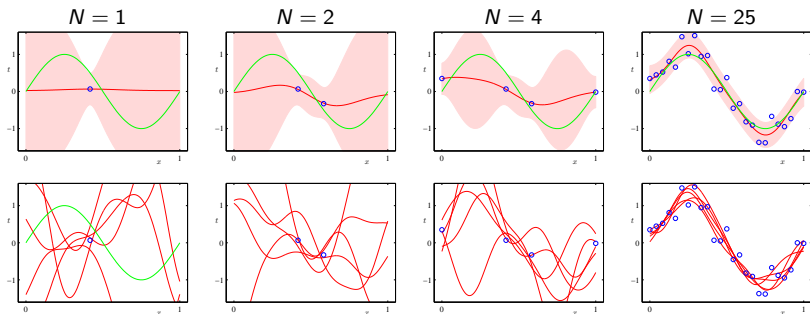
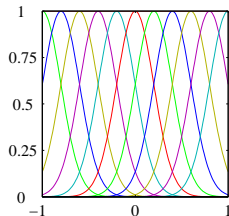
$$p(y^*|\mathbf{y}) = \int p(y^*|\mathbf{x}^*, \mathbf{w})p(\mathbf{w}|\mathbf{y})d\mathbf{w} \quad (\textit{Posterior predictive dist.})$$

- Spend 5 minutes DTU Learn quiz: "Lecture 3: Bayesian inference"

## Example: posterior predictive distributions

- *Predictive distributions* for simple sinoidal toy dataset using Gaussian Basis functions
- Samples from the posterior

Feature extractors  $\phi_i(x)$



## Dealing with hyperparameters

## But what about the hyperparameters?

- The Bayesian linear regression model

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) \quad (\text{prior})$$

$$p(\mathbf{y}|\mathbf{w}) = \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2\mathbf{I}) \quad (\text{likelihood})$$

$$p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S}) \quad (\text{posterior})$$

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \sigma^2\mathbf{I} + \alpha^{-1}\Phi\Phi^T) \quad (\text{marginal likelihood})$$

- A fully Bayesian solution would require us to impose priors on  $\alpha, \beta$  (recall  $\beta = \frac{1}{\sigma^2}$ )

$$p(\alpha, \beta|\mathbf{y}) \propto p(\mathbf{y}|\alpha, \beta)p(\alpha, \beta)$$

- ... and integrate them out

$$\begin{aligned} p(y^*|\mathbf{y}) &= \iiint p(y^*|\mathbf{w}, \beta)p(\mathbf{w}|\mathbf{y}, \alpha, \beta)p(\alpha, \beta|\mathbf{y}) \, d\mathbf{w}d\alpha d\beta \\ &= \mathbb{E}_{p(\alpha, \beta|\mathbf{y})} \left[ \mathbb{E}_{p(\mathbf{w}|\mathbf{y}, \alpha, \beta)} [p(y^*|\mathbf{w}, \beta)] \right] \end{aligned}$$

- ... but this is not analytically tractable. Later in the course we will learn tools to deal with this in general

## The evidence approximation

- How to deal with this bastard?

$$\begin{aligned} p(y^*|\mathbf{y}) &= \int \int \int p(y^*|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{y}, \alpha, \beta) p(\alpha, \beta|\mathbf{y}) d\mathbf{w} d\alpha d\beta \\ &= \mathbb{E}_{p(\alpha, \beta|\mathbf{y})} \left[ \mathbb{E}_{p(\mathbf{w}|\mathbf{y}, \alpha, \beta)} [p(y^*|\mathbf{w}, \beta)] \right] \end{aligned}$$

- *The evidence approximation, maximum likelihood type II, Empirical Bayes*
- If the posterior is sharply peaked around  $\hat{\alpha}$  and  $\hat{\beta}$ , then  $p(\alpha, \beta|\mathbf{y})$  can be approximated by a Dirac's delta distribution

$$p(y^*|\mathbf{y}) \approx \mathbb{E}_{p(\mathbf{w}|\mathbf{y}, \hat{\alpha}, \hat{\beta})} [p(y^*|\mathbf{w}, \beta)]$$

- Assume we impose a *flat prior* on  $\alpha$  and  $\beta$ , then

$$p(\alpha, \beta|\mathbf{y}) \propto p(\mathbf{y}|\alpha, \beta) p(\alpha, \beta) \propto p(\mathbf{y}|\alpha, \beta)$$

- We can estimate  $\hat{\alpha}, \hat{\beta}$  by *optimizing the marginal likelihood*  $p(\mathbf{y}|\alpha, \beta)$

$$\hat{\alpha}, \hat{\beta} = \arg \max_{\alpha, \beta} \log p(\mathbf{y}|\alpha, \beta)$$

## Sinoidal example revisited using the evidence approximation

