

# Bayesian Learning

Lecture series „Machine Learning“

Niels Landwehr

Research Group „Data Science“  
Institute of Computer Science  
University of Hildesheim

# Agenda for Lecture

- Fundamental concepts of Bayesian learning
- Introductory example: coin tosses
- Bayesian linear regression

# Agenda for Lecture

- Fundamental concepts of Bayesian learning
- Introductory example: coin tosses
- Bayesian linear regression

# Review: Probabilistic Model for Data

- We are still concerned with supervised learning problems:
  - Models  $f : \mathcal{X} \rightarrow \mathcal{Y}$
  - Training data  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$
- Review: probabilistic assumptions about training data
  - Training instances independently drawn from joint distribution over inputs and outputs:

$$(\mathbf{x}_n, y_n) \sim p(\mathbf{x}, y)$$

- According to product rule, can split up joint distribution into

- The instances  $\mathbf{x}_n$  are sampled from a probability distribution over instances.
- $p(\mathbf{x})$  describes distribution over population of objects

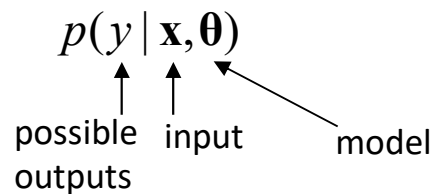
$$\mathbf{x}_n \sim p(\mathbf{x})$$

$$y_n \sim p(y | \mathbf{x}_n)$$

- Given an instance  $\mathbf{x}_n$ , its label is drawn from a distribution  $p(y | \mathbf{x}_n)$  that represents the relationship between input and output.

# Review: Probabilistic Models

- Review: probabilistic models such as logistic regression define a conditional probability distribution

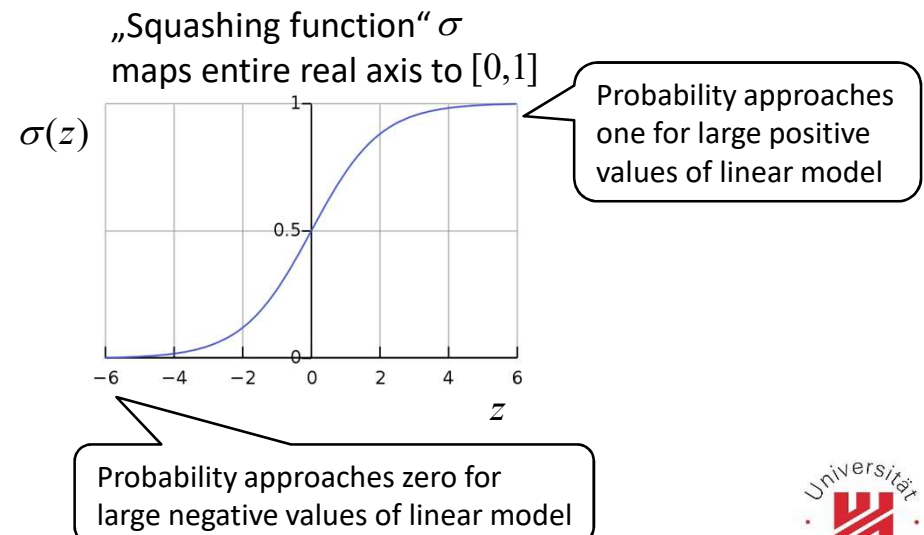


- For example, binary logistic regression: distribution over binary  $y \in \{0,1\}$  by defining probability for positive class as sigmoid function of linear model

$$p(y = 1 | \mathbf{x}, \boldsymbol{\theta}) = \sigma(\mathbf{x}^T \boldsymbol{\theta})$$

$$p(y = 0 | \mathbf{x}, \boldsymbol{\theta}) = 1 - \sigma(\mathbf{x}^T \boldsymbol{\theta})$$

$$\sigma(z) = \frac{e^z}{1 + e^z}$$



# Probabilistic Models: Maximum Likelihood

- So far, we have learned probabilistic models by maximum likelihood (lecture on linear classification):

**Optimization Problem:**

$$\begin{aligned}\boldsymbol{\theta}^* &= \arg \max_{\boldsymbol{\theta}} p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \prod_{n=1}^N p(y_n \mid \mathbf{x}_n, \boldsymbol{\theta})\end{aligned}$$

- This yields a single final model (parameter estimate)  $\boldsymbol{\theta}^*$
- Question: can we really be sure that  $\boldsymbol{\theta}^*$  is the best model?
  - The selection of the best model  $\boldsymbol{\theta}^*$  out of all possible model parameters is based on a limited set of training data  $\mathcal{D}$
  - Realistically, there is remaining uncertainty: even after having seen the training data, cannot be 100% sure what the best model is
  - **Idea: can we quantify this remaining uncertainty?**

# Bayesian Model of Data Generation

- **Bayesian Learning:** extend the probabilistic model of how the training (and test) data are generated
- Specifically, reason probabilistically about models as well, not only data
- Statistically speaking, the model itself becomes a random variable
- **Bayesian model of how the data (and the true model) is generated:**
  1. The true model  $\theta^*$  is drawn from a so-called **prior distribution**  $p(\theta)$  over possible models.
    - The true model  $\theta^*$  is unknown, but we make an assumption for the prior distribution  $p(\theta)$
    - Assumption could be, for example, that models with small parameter values are more likely than models with large parameter values (idea of shrinkage)
  2. The inputs  $\mathbf{x}_1, \dots, \mathbf{x}_N$  are drawn independently from a distribution over inputs:
  3. The outputs  $y_1, \dots, y_N$  are then drawn from the distribution given by the model  $\theta^*$

$$\mathbf{x}_n \sim p(\mathbf{x})$$

$$y_n \sim p(y | \mathbf{x}, \theta^*)$$

# Bayes Rule and Posterior Probability of Model

- At the heart of Bayesian machine learning is the **Bayes rule**, a simple rule for any two random variables  $u, v$  which have a joint distribution  $p(u, v)$ :

$$p(u | v) = \frac{p(v | u)p(u)}{p(v)}$$

- In Bayesian learning, we apply this rule to models and data in the following way (rough idea, see below for a more formal/detailed treatment):

Bayes rule computes the so-called **posterior probability** of a model given the data: How likely is it that a particular model is the true model, given the data we have seen? The distribution tells us both about likely correct models and the remaining uncertainty.

This is the **likelihood**: which probability does the model assign to the training data? Can be computed given model and data.

This is the **prior**: what are a priori likely models (before seeing the data). Assumed known.

$$p(\text{Model} | \text{Data}) = \frac{p(\text{Data} | \text{Model})p(\text{Model})}{p(\text{Data})}$$

Normalizing factor that can in principle be computed from likelihood and prior



# Bayes Rule and Posterior Probability of Model

- More formally, for a model of the form  $p(y | \mathbf{x}, \boldsymbol{\theta})$ , the Bayes rule reads as follows:

The diagram shows the equation  $p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X}) p(\boldsymbol{\theta})}{p(\mathbf{y} | \mathbf{X})}$  with four callout boxes:

- Top-left: Models are given by parameter vectors (points to  $\boldsymbol{\theta}$ )
- Top-middle: Likelihood: probability of observed labels given input and model (points to  $p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X})$ )
- Top-right: Prior over model parameters (points to  $p(\boldsymbol{\theta})$ )
- Bottom-right: Probability of data, only normalizing factor that is independent of model  $\boldsymbol{\theta}$  (points to  $p(\mathbf{y} | \mathbf{X})$ )

- Here,  $\mathbf{X}$  and  $\mathbf{y}$  are the training data in the usual matrix/vector notation
- Proof: simply Bayes rule applied to the variables  $\boldsymbol{\theta}, \mathbf{y}$  with joint distribution given by the conditional  $p(\boldsymbol{\theta}, \mathbf{y} | \mathbf{X})$  (note that  $p(\boldsymbol{\theta} | \mathbf{X}) = p(\boldsymbol{\theta})$  because prior is independent of input data  $\mathbf{X}$ )
- Note that because the model only describes a distribution  $p(y | \mathbf{x}, \boldsymbol{\theta})$  over outputs given inputs and parameters, „Data“ in this case only refers to the observed labels  $\mathbf{y}$

# Bayesian Predictions

- Bayes rule gives us a posterior probability distribution over models given the data
- How do we obtain predictions for a novel test instance  $\mathbf{x}_{new}$ ?
- **First approach:**
  - Find the most probable model given the prior and the data:

$$\begin{aligned}\boldsymbol{\theta}_{MAP} &= \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) \\ &= \arg \max_{\boldsymbol{\theta}} \frac{p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{y} | \mathbf{X})} \\ &= \arg \max_{\boldsymbol{\theta}} p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta})\end{aligned}$$

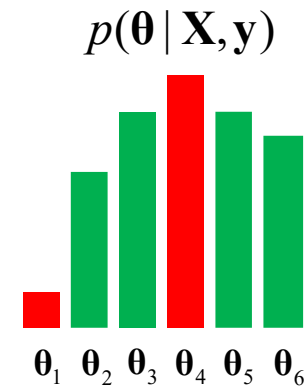
The model  $\boldsymbol{\theta}_{MAP}$  is also called the **maximum a posteriori** (MAP) model

- We can then compute a prediction  $\hat{y}$  using this model:

$$\hat{y} = \arg \max_y p(y | \mathbf{x}_{new}, \boldsymbol{\theta}_{MAP})$$

# Motivation: Bayesian Predictions

- The prediction with the maximum-a-posteriori model does not take into account the remaining uncertainty inherent in the posterior distribution
- Let's look at a toy example:
  - Assume that our model space only contains six models overall:  $\{\theta_1, \theta_2, \theta_3, \theta_4, \theta_5, \theta_6\}$
  - Of those six models, two predict the negative class (red) and four predict the positive class (green)
  - Model  $\theta_4$  is the maximum a posteriori model, and would result in the prediction „negative class“
  - However, is this the right decision? There is a lot of remaining uncertainty, and most other models would predict the positive class...
- Would prefer to better take remaining uncertainty into account
- **Idea:** take into account predictions of **all models**, and weight them by their posterior probability (see next slide)



# Bayesian Predictions

- **Second approach: Bayesian prediction**
- Idea: we do not limit ourselves to one model, but directly compute the most probable prediction given the evidence from the training data and the new instance:

Most probable prediction, given training data  $\mathbf{X}, \mathbf{y}$  and the new instance  $\mathbf{x}_{new}$

$$\hat{y} = \arg \max_y p(y | \mathbf{x}_{new}, \mathbf{X}, \mathbf{y})$$

Marginalization rule  $p(a) = \int p(a, b) db$   
applied to conditional distribution  $p(y, \boldsymbol{\theta} | \mathbf{x}_{new}, \mathbf{X}, \mathbf{y})$

$$= \arg \max_y \int p(y, \boldsymbol{\theta} | \mathbf{x}_{new}, \mathbf{X}, \mathbf{y}) d\boldsymbol{\theta}$$

Product rule  $p(a, b) = p(a | b)p(b)$  applied to  
conditional distribution  $p(y, \boldsymbol{\theta} | \mathbf{x}_{new}, \mathbf{X}, \mathbf{y})$

$$= \arg \max_y \int p(y | \boldsymbol{\theta}, \mathbf{x}_{new}, \mathbf{X}, \mathbf{y}) p(\boldsymbol{\theta} | \mathbf{x}_{new}, \mathbf{X}, \mathbf{y}) d\boldsymbol{\theta}$$

$$= \arg \max_y \int p(y | \mathbf{x}_{new}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) d\boldsymbol{\theta}$$

Given the model, the training data does not influence the prediction:

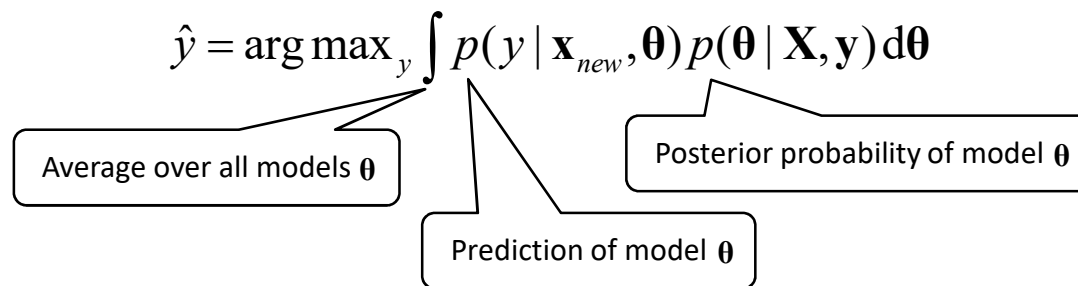
$$p(y | \boldsymbol{\theta}, \mathbf{x}_{new}, \mathbf{X}, \mathbf{y}) = p(y | \boldsymbol{\theta}, \mathbf{x}_{new})$$

The learned model is independent of the new test input:

$$p(\boldsymbol{\theta} | \mathbf{x}_{new}, \mathbf{X}, \mathbf{y}) = p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y})$$

# Bayesian Predictions

- **Bayesian prediction: average over the predictions of all models, weighted by their posterior probability**

$$\hat{y} = \arg \max_y \int p(y | \mathbf{x}_{new}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) d\boldsymbol{\theta}$$


Average over all models  $\boldsymbol{\theta}$

Posterior probability of model  $\boldsymbol{\theta}$

Prediction of model  $\boldsymbol{\theta}$

- Bayesian prediction takes into account the predictive distributions of all models
- Models are weighted by their posterior probability

# Agenda for Lecture

- Fundamental concepts of Bayesian learning
- **Introductory example: coin tosses**
- Bayesian linear regression

# Example: Coin Tosses

- Introductory example for Bayesian learning: Tossing a binary coin  $N$  times
- Coin tossing experiment
  - Assume coin that can be tossed and then lands on head ( $y = 1$ ) or tail ( $y = 0$ )
  - The coin is not necessarily fair, and the probability of landing on head is  $\theta \in [0,1]$ :

$$p(y = 1 | \theta) = \theta$$

$$p(y = 0 | \theta) = 1 - \theta$$

- Assume we are tossing the coin  $N$  times, observing results  $\mathbf{y} = (y_1, \dots, y_N)^T$  with  $y_i \in \{0,1\}$
- We can summarize the observations into counts for the outcomes head and tail:

$$N_h = \sum_{n=1}^N I(y_n = 1) \quad \text{number of heads}$$

$$N_t = \sum_{n=1}^N I(y_n = 0) \quad \text{number of tails}$$

# Coin Tosses: Bayesian Parameter Estimates

- Coin tosses can be seen as a simple machine learning setting:
  - Unknown parameter  $\theta$  is the model
  - Head and tail counts  $N_h, N_t$  are the training observations
  - Task is to learn model  $\theta$  from training observations
- We want to follow a Bayesian approach:
  - What can we learn from the data about the model parameter  $\theta$  ?
  - The data is given by the head and tail counts, and knowing  $N_h$  is sufficient given a fixed  $N$
  - Bayes rule:

$$p(\theta | N_h) = \frac{p(N_h | \theta) p(\theta)}{p(N_h)}$$

Posterior probability of a possible coin parameter  $\theta$  after having seen the data

Likelihood: probability of observed outcomes given a parameter  $\theta$

We need a prior over possible coin parameters  $\theta$

Probability of data, only normalizing factor that is independent of model  $\theta$



# Coin Tosses: Likelihood

- Likelihood  $p(N_h | \theta)$  : what is the probability of seeing exactly  $N_h$  heads (and therefore  $N_t = N - N_h$  tails) for a given coin parameter  $\theta$  ?
- Likelihood is given by the Binomial distribution:

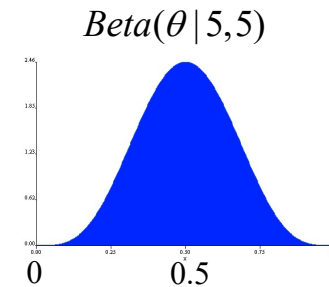
$$\text{Bin}(N_h | N, \theta) = \underbrace{\binom{N}{N_h}}_{\text{number of series with } N_h \text{ heads}} \underbrace{\theta^{N_h} (1 - \theta)^{N - N_h}}_{\text{probability for a specific series of tosses with } N_h \text{ heads and } N_t = N - N_h \text{ tails}}$$

- "Binomial coefficient"  $\binom{N}{N_h} = \frac{N!}{N_h! (N - N_h)!}$

# Coin Tosses: Prior

- Prior  $p(\theta)$ : what is a good choice for a prior distribution over coin toss parameters  $\theta$ ?
- Good choice for prior distribution is the so-called Beta distribution:

$$\begin{aligned} p(\theta) &= \text{Beta}(\theta | \alpha_h, \alpha_t) \\ &= \frac{\Gamma(\alpha_h + \alpha_t)}{\Gamma(\alpha_h)\Gamma(\alpha_t)} \theta^{\alpha_h-1} (1-\theta)^{\alpha_t-1} \end{aligned}$$



- Here,  $\Gamma : \mathbb{R} \rightarrow \mathbb{R}$  is a continuous extension of the factorial function:

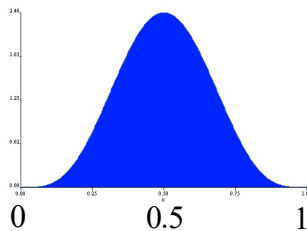
$$\Gamma(z) = \int_0^{\infty} t^{z-1} e^{-t} dt \quad \forall n \in \mathbb{N} : \Gamma(n) = (n-1)!$$

- The values  $\alpha_h$ ,  $\alpha_t$  are parameters of the Beta distribution that determine the shape of the prior. Called hyperparameters to separate them from model parameter  $\theta$

# Coin Tosses: Prior

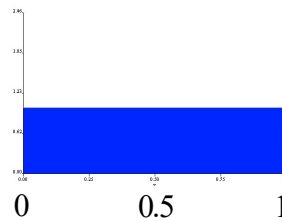
- Examples for Beta distributions with different hyperparameters:

$$\alpha_h = 5, \alpha_t = 5$$



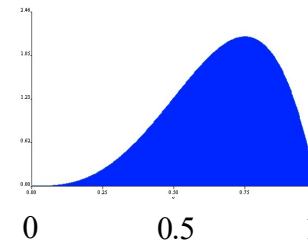
Prior assumption:  
Most probably the coin is approximately fair, unlikely to always show head or always show tail

$$\alpha_h = 1, \alpha_t = 1$$



Prior assumption:  
No prior assumptions about parameter, so-called uninformative prior

$$\alpha_h = 4, \alpha_t = 2$$



Prior assumption:  
Coin will probably be biased towards landing head

# Why a Beta Prior?

- Why did we choose this prior distribution?
- Structural similarity to likelihood function:

**Prior:**  $p(\theta) = \text{Beta}(\theta | \alpha_h, \alpha_t) = \frac{\Gamma(\alpha_h + \alpha_t)}{\Gamma(\alpha_h)\Gamma(\alpha_t)} \theta^{\alpha_h-1} (1-\theta)^{\alpha_t-1}$

**Likelihood:**  $p(N_h | \theta) = \text{Bin}(N_h | N, \theta) = \binom{N}{N_h} \theta^{N_h} (1-\theta)^{N_t}$

- This makes it easier to compute the posterior distribution over parameters given the empirical data (number of heads): according to Bayes rule,

$$p(\theta | N_h) = \frac{1}{Z} p(N_h | \theta) p(\theta)$$

Posterior probability of a possible coin parameter  $\theta$  after having seen the data

Normalizer

Likelihood: probability of observed outcomes given a parameter  $\theta$

Prior over possible coin parameters  $\theta$

# Posterior is Again Beta Distribution

- Let's compute the posterior using Bayes rule:

$$\begin{aligned} p(\theta | N_h) &= \frac{p(N_h | \theta)p(\theta)}{p(N_h)} \\ &= \frac{1}{Z} \text{Bin}(N_h | N, \theta) \text{Beta}(\theta | \alpha_h, \alpha_t) \\ &= \frac{1}{Z} \binom{N}{N_h} \theta^{N_h} (1 - \theta)^{N_t} \frac{\Gamma(\alpha_h + \alpha_t)}{\Gamma(\alpha_h)\Gamma(\alpha_t)} \theta^{\alpha_h - 1} (1 - \theta)^{\alpha_t - 1} \\ &= \frac{1}{Z'} \theta^{\alpha_h + N_h - 1} (1 - \theta)^{\alpha_t + N_t - 1} \end{aligned}$$

What is  $Z'$ ?  
What is the correct normalizer such that the final distribution integrates to one?

# Posterior is Again Beta Distribution

- Let's compute the posterior using Bayes rule:

$$\begin{aligned} p(\theta | N_h) &= \frac{p(N_h | \theta) p(\theta)}{p(N_h)} \\ &= \frac{1}{Z} \text{Bin}(N_h | N, \theta) \text{Beta}(\theta | \alpha_h, \alpha_t) \\ &= \frac{1}{Z} \binom{N}{N_h} \theta^{N_h} (1 - \theta)^{N_t} \frac{\Gamma(\alpha_h + \alpha_t)}{\Gamma(\alpha_h) \Gamma(\alpha_t)} \theta^{\alpha_h - 1} (1 - \theta)^{\alpha_t - 1} \\ &= \frac{1}{Z'} \theta^{\alpha_h + N_h - 1} (1 - \theta)^{\alpha_t + N_t - 1} \\ &= \frac{\Gamma(\alpha_h + N_h + \alpha_t + N_t)}{\Gamma(\alpha_h + N_h) \Gamma(\alpha_t + N_t)} \theta^{\alpha_h + N_h - 1} (1 - \theta)^{\alpha_t + N_t - 1} \\ &= \text{Beta}(\theta | \alpha_h + N_h, \alpha_t + N_t) \end{aligned}$$

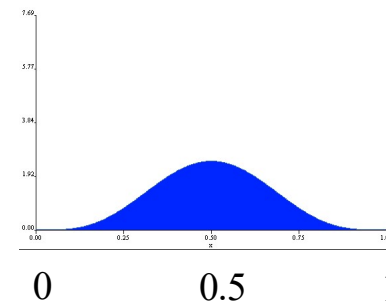
Has to be this normalizer, because this is the normalizer of the Beta distribution with hyperparameters  $\alpha_h + N_h$ ,  $\alpha_t + N_t$

- Posterior is again a Beta distribution, with new hyperparameters  $\alpha_h + N_h$ ,  $\alpha_t + N_t$
- We call this a conjugate prior: posterior is in the same distribution family as prior**

# Example: Posterior For Coin Tosses

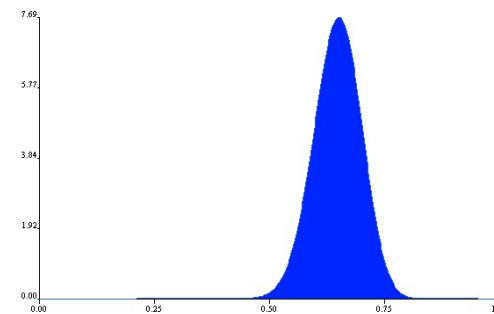
- Example: Computing posterior for coin tosses
- Let's say we choose a prior distribution  $p(\theta) = \text{Beta}(\theta | 5, 5)$

Prior expresses the belief that fair coins are more probable than unfair coins, and very unfair coins ( $\theta$  near zero or one) are very improbable



- Let's assume we have thrown the coin  $N = 75$  times, and have observed  $N_h = 50$  times „Head“ and  $N_t = 25$  times „Tail“
- According to posterior calculation on last slide, the posterior is again a Beta distribution with new hyperparameters  $5+50, 5+25$ :  $\text{Beta}(\theta | 55, 30)$

Posterior belief: coin is probably unfair ( $\theta > 0.5$ ).  
There is some remaining uncertainty about how unfair the coin is



# Coin Tosses: Bayesian Prediction

- Bayesian prediction for coin tosses:
  - Assume we have seen the result of  $N$  coin tosses in the training data
  - What is the probability for a new toss of the same coin to land on „Head“ or „Tail“?
- Bayesian approach: do not use a single model to predict, but average over models weighted by posterior

$N$  observed tosses (training data)

New coin toss

$$\begin{aligned} p(y_{new} = 1 | \mathbf{y}) &= \int p(y_{new} = 1 | \theta) p(\theta | \mathbf{y}) d\theta \\ &= \int \theta \text{Beta}(\theta | \alpha_h + N_h, \alpha_t + N_t) d\theta && p(y_{new} = 1 | \theta) = \theta, \text{ plug in posterior} \\ &= \frac{\alpha_h + N_h}{\alpha_h + \alpha_t + N_h + N_t} && \begin{array}{l} \text{Expectation of Beta distribution with} \\ \text{hyperparameters } \alpha_h + N_h, \alpha_t + N_t. \\ \text{Standard result, no proof here.} \end{array} \end{aligned}$$

- Bayesian prediction for the example of last slide: probability for next coin toss to land on head would be  $55 / 85 \approx 0.647$



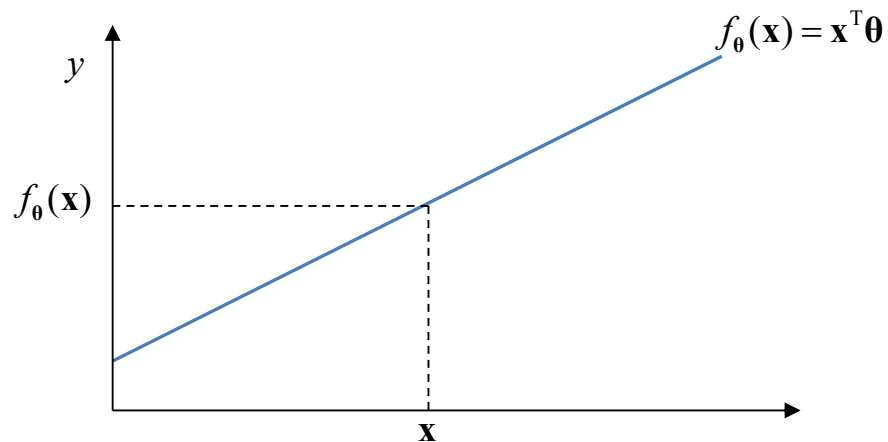
# Agenda for Lecture

- Fundamental concepts of Bayesian learning
- Introductory example: coin tosses
- Bayesian linear regression

# Review: Linear Regression

- Review: linear regression model
  - Defines model  $f_{\theta} : \mathbb{R}^M \rightarrow \mathbb{R}$
  - Prediction is linear function of input attributes:

$$\begin{aligned} f_{\theta}(\mathbf{x}) &= \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_M x_M \\ &= \mathbf{x}^T \boldsymbol{\theta} \end{aligned}$$



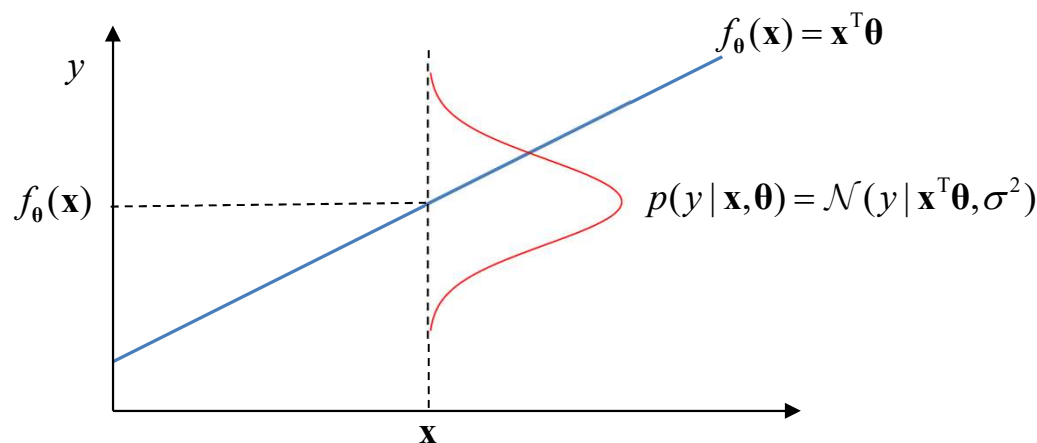
# Linear Regression as a Probabilistic Model

- Linear regression can be extended to a probabilistic model as follows
- Define a distribution over the output  $y \in \mathbb{R}$  given the input  $\mathbf{x} \in \mathbb{R}^M$  :

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

Normal distribution with mean  $\mathbf{x}^T \boldsymbol{\theta}$  and variance  $\sigma^2$ .  
Variance  $\sigma^2$  is a hyperparameter of model

- Encodes some uncertainty about prediction



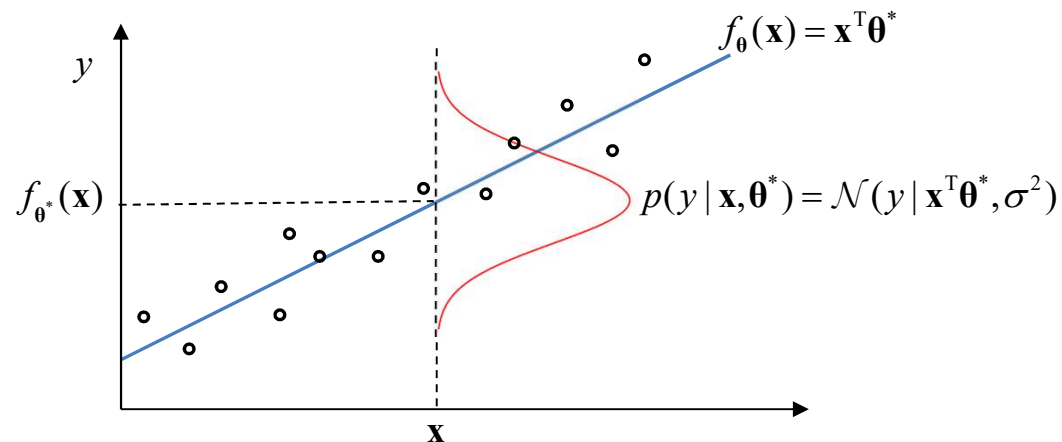
# Linear Regression as a Probabilistic Model

- Remember that in a Bayesian setting, we assume that the data (labels) have been generated from an unknown true model  $\theta^*$ :

$$y_n \sim p(y | \mathbf{x}, \theta^*)$$

- In this case, this means that we assume that the labels are generated from a linear model plus Gaussian noise:

$$y_n = \mathbf{x}_n^T \theta^* + \epsilon_n \quad \epsilon_n \sim \mathcal{N}(\epsilon | 0, \sigma^2)$$



# Bayes Rule for Probabilistic Linear Regression

- We use Bayes rule to compute the posterior probability:

Posterior distribution over model parameters

Likelihood: probability of observed labels given inputs and model

We need a prior over model parameters

$$p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X}) p(\boldsymbol{\theta})}{p(\mathbf{y} | \mathbf{X})}$$

Probability of data, only normalizing factor that is independent of model  $\boldsymbol{\theta}$

- We need to
  - compute likelihood
  - define prior distribution
  - derive posterior distribution

# Likelihood for Probabilistic Linear Regression

- Likelihood for probabilistic regression model:

$$p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \boldsymbol{\theta})$$

$$= \prod_{n=1}^N p(y_n | \mathbf{x}_n, \boldsymbol{\theta})$$

Independent training instances

$$= \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^T \boldsymbol{\theta}, \sigma^2)$$

Plugging in model

$$= \mathcal{N}(\mathbf{y} | \mathbf{X}^T \boldsymbol{\theta}, \sigma^2 \mathbf{I})$$

Product of univariate normal distributions can be written as multivariate normal distribution:

- Mean vector is vector of predictions:

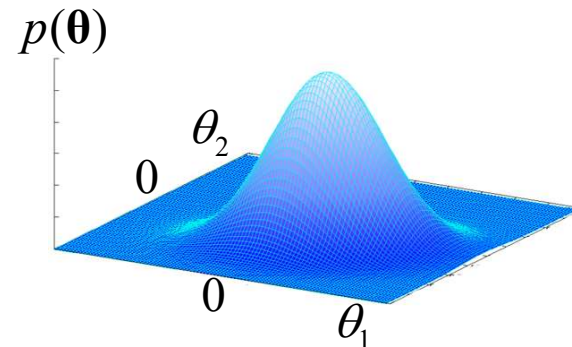
$$\mathbf{X}^T \boldsymbol{\theta} = \begin{pmatrix} \mathbf{x}_1^T \boldsymbol{\theta} \\ \dots \\ \mathbf{x}_N^T \boldsymbol{\theta} \end{pmatrix}$$

- Covariance matrix is a diagonal matrix  $\sigma^2 \mathbf{I}$

# Prior for Probabilistic Linear Regression

- Prior for probabilistic regression model:
  - We need to define a prior over parameter vectors  $\boldsymbol{\theta} \in \mathbb{R}^M$
  - Similar idea as for regularization/shrinkage: we prefer parameter values that are not too large (in absolute value)
- As prior distribution, choose a multivariate normal distribution with mean zero and a diagonal covariance matrix:

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, \sigma_p^2 \mathbf{I})$$



- Hyperparameter  $\sigma_p^2$  controls variance of prior distribution and thereby strength of preference for small parameters
  - small  $\sigma_p^2$  : low variance, very „peaked“ distribution, strong regularization
  - large  $\sigma_p^2$  : high variance, flat distribution, less regularization

# Posterior for Probabilistic Linear Regression

- The prior distribution is a conjugate prior for the likelihood computed above: the product of prior and likelihood is again a multivariate normal distribution

$$\begin{aligned} p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) &= \frac{1}{Z} p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X}) p(\boldsymbol{\theta}) && \text{Bayes rule} \\ &= \frac{1}{Z} \mathcal{N}(\mathbf{y} | \mathbf{X}^T \boldsymbol{\theta}, \sigma^2 \mathbf{I}) \cdot \mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, \sigma_p^2 \mathbf{I}) && \text{plugging in likelihood and prior} \\ &= \mathcal{N}(\boldsymbol{\theta} | \bar{\boldsymbol{\theta}}, \mathbf{A}^{-1}) && \text{compute product of normal terms - no details here} \end{aligned}$$

$$\text{where } \mathbf{A} = \sigma^{-2} \mathbf{X}^T \mathbf{X} + \sigma_p^{-2} \mathbf{I} \quad \text{and} \quad \bar{\boldsymbol{\theta}} = \sigma^{-2} \mathbf{A}^{-1} \mathbf{X}^T \mathbf{y}$$

- Posterior is again normally distributed, with a new mean vector given by  $\bar{\boldsymbol{\theta}}$  and a new covariance matrix  $\mathbf{A}^{-1}$



# Posterior for Probabilistic Linear Regression

- The maximum a posteriori model is

$$\begin{aligned}\boldsymbol{\theta}_{MAP} &= \arg \max_{\boldsymbol{\theta}} \mathcal{N}(\boldsymbol{\theta} | \bar{\boldsymbol{\theta}}, \mathbf{A}^{-1}) \\ &= \bar{\boldsymbol{\theta}}\end{aligned}$$

- Note the similarity to the solution of ridge regression:
  - MAP solution:

$$\boldsymbol{\theta}_{MAP} = \sigma^{-2} (\sigma^{-2} \mathbf{X}^T \mathbf{X} + \sigma_p^{-2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

- Ridge regression (squared loss, L2-regularization):

$$\boldsymbol{\theta}^* = (\mathbf{X}^T \mathbf{X} + N\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

- For appropriately chosen hyperparameters, the solution is the same