

Photogrammetry & Robotics Lab

Machine Learning for Robotics and Computer Vision

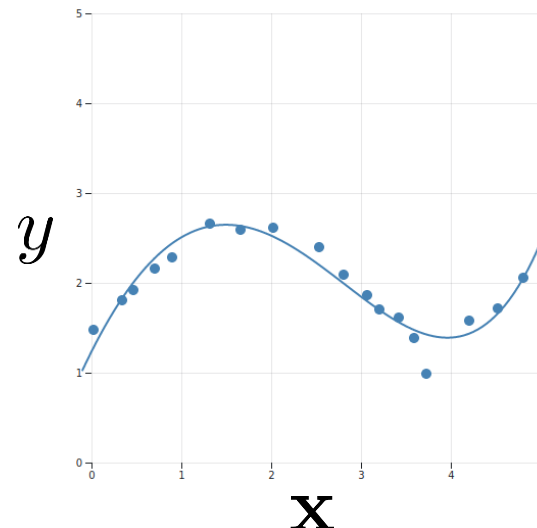
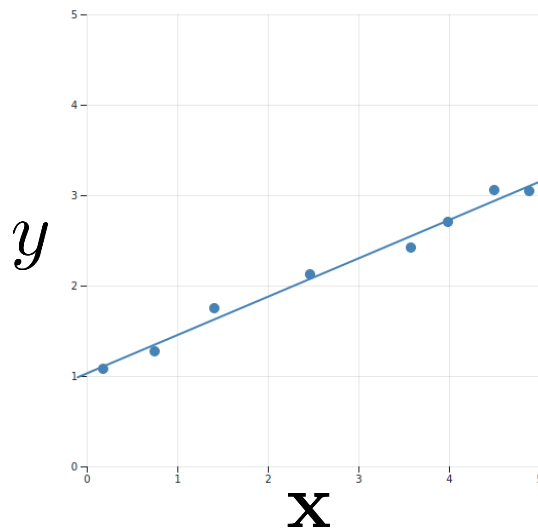
Regression

Jens Behley

Recap: Last Lecture

- High-level overview of machine learning algorithms
- Main ingredients:
 1. **Data**
 2. **Model**
 3. **Learning**
- Discussed the importance of train, validation, and test set
- Discussed as an example **k-Nearest Neighbor** classification

Regression



- **Regression** is finding a function $f(\mathbf{x})$ that explains our targets $y \in \mathbb{R}$ for an input $\mathbf{x} \in \mathbb{R}^D$
- Assumption: we have noisy observations:

$$y_n = f(\mathbf{x}_n) + \epsilon \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- We assume that we know σ^2 in advance.

Linear Regression

- Under this assumptions, this leads to the following probabilistic formulation

$$P(y|\mathbf{x}, \theta) = \mathcal{N}(y|f(\mathbf{x}), \sigma^2)$$

- In linear regression, we assume that parameters θ appear **linearly** in our model

$$f(\mathbf{x}) = \mathbf{x}^T \theta + \theta_0$$

- θ_0 is called **intercept** (or **bias**) that enables us to have also functions that do not pass through the origin

Some Notation

- Training data given by

$$\mathcal{X}_{\text{train}} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n), \dots, (\mathbf{x}_N, y_N)\}$$

with $\mathbf{x}_n = (x_1, \dots, x_d, \dots, x_D) \in \mathbb{R}^D$ and $y \in \mathbb{R}$

- To simplify notation, we will add $x_0 = 1$
- Which simplifies $f(\mathbf{x})$ as follows:

$$f(\mathbf{x}) = \mathbf{x}^T \theta + \theta_0 \xrightarrow{\mathbf{x} := (1, \mathbf{x})^T} f(\mathbf{x}) = \mathbf{x}^T \theta$$

- Define $\mathbf{x}_{1:N} := \mathbf{x}_1, \dots, \mathbf{x}_N$ and $y_{1:N} := y_1, \dots, y_N$

Probabilistic view

- Treat parameters θ as random variables
- We are interested in the posterior

$$P(\theta|\mathbf{x}_{1:N}, y_{1:N}) = \frac{P(y_{1:N}|\mathbf{x}_{1:N}, \theta)P(\theta)}{P(y_{1:N}|\mathbf{x}_{1:N})}$$

- As $P(y_{1:N}|\mathbf{x}_{1:N})$ is independent of θ , follows:

$$\underbrace{P(\theta|\mathbf{x}_{1:N}, y_{1:N})}_{\text{Posterior}} \propto \underbrace{P(y_{1:N}|\mathbf{x}_{1:N}, \theta)}_{\text{Likelihood}} \underbrace{P(\theta)}_{\text{Prior}}$$

- $P(y_{1:N}|\mathbf{x}_{1:N})$ is marginal likelihood/evidence

Common Parameter Estimation

- Learning is finding parameters of

$$P(\theta|\mathbf{x}_{1:N}, y_{1:N}) \propto P(y_{1:N}|\mathbf{x}_{1:N}, \theta) \quad P(\theta)$$

- Paradigms for parameter estimation:
 1. Point estimate with uniform prior
→ Maximum Likelihood Estimation
 2. Point estimate with given prior
→ Maximum A posteriori Estimation (MAP)
 3. Determine posterior over the parameters
→ Bayesian Estimation

Maximum Likelihood Estimation

- Assuming a uniform prior, the posterior reduces to

$$P(\theta|\mathbf{x}_{1:N}, y_{1:N}) \propto P(y_{1:N}|\mathbf{x}_{1:N}, \theta)$$

- We want to find the parameters θ^* that maximize the likelihood

$$\begin{aligned}\theta^* &= \arg \max_{\theta} P(y_{1:N}|\mathbf{x}_{1:N}, \theta) \\ &= \arg \max_{\theta} \prod_{n=1}^N P(y_n|\mathbf{x}_n, \theta) \quad (\text{i.i.d.})\end{aligned}$$

Negative Log-Likelihood

- As logarithm is a monotonically increasing function \rightarrow optimum of $f \Leftrightarrow$ optimum of $\log f$
- With the negative log-transform, we now **minimize** accordingly:

$$\begin{aligned}\theta^* &= \arg \min_{\theta} - \log \prod_{n=1}^N P(y_n | \mathbf{x}_n, \theta) \\ &= \arg \min_{\theta} - \underbrace{\sum_{n=1}^N \log P(y_n | \mathbf{x}_n, \theta)}\end{aligned}$$

Negative Log-Likelihood

- Advantage: sum is numerical more stable than product of values in $[0,1]$!

NLL for Linear Regression

- Inserting all terms for Linear Regression:

$$\begin{aligned} P(y_n | \mathbf{x}_n, \theta) &= \mathcal{N}(y_n | f(\mathbf{x}_n), \sigma^2) \\ &= \frac{1}{\sigma \sqrt{2\pi}} \exp \left(-\frac{1}{2} \frac{(y_n - f(\mathbf{x}_n))^2}{\sigma^2} \right) \\ &= \frac{1}{\sigma \sqrt{2\pi}} \exp \left(-\frac{1}{2} \frac{(y_n - \mathbf{x}_n^T \theta)^2}{\sigma^2} \right) \quad (f(\mathbf{x}_n) = \mathbf{x}_n^T \theta) \end{aligned}$$

- It follows:

$$\log P(y_n | \mathbf{x}_n, \theta) = \text{const} - \frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^T \theta)^2$$

- The NLL $\mathcal{L}(\theta)$ is then:

$$\mathcal{L}(\theta) := \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^T \theta)^2$$

Design matrix

- Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ and $\mathbf{y} \in \mathbb{R}^N$ be defined as:

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \in \mathbb{R}^{N \times D} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \in \mathbb{R}^N$$

- $\mathbf{X} \in \mathbb{R}^{N \times D}$ is called the **design matrix**
- Using these definitions, we can rewrite:

$$\begin{aligned} \mathcal{L}(\theta) &:= \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^T \theta)^2 \\ &= \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\theta)^T (\mathbf{y} - \mathbf{X}\theta) \end{aligned}$$

Gradient of NLL

- Minimizing the NLL $\mathcal{L}(\theta)$ means that we have to find θ where gradient $\frac{d\mathcal{L}}{d\theta}$ is zero
- Gradient can be derived as follows:

$$\begin{aligned}\frac{d\mathcal{L}}{d\theta} &= \frac{d}{d\theta} \left(\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\theta)^T (\mathbf{y} - \mathbf{X}\theta) \right) \\ &= \frac{1}{2\sigma^2} \frac{d}{d\theta} \left(\cancel{\mathbf{y}^T \mathbf{y}} - \underline{2\mathbf{y}^T \mathbf{X}\theta} + \underline{\theta^T \mathbf{X}^T \mathbf{X}\theta} \right) \\ &= \frac{1}{2\sigma^2} \left(\underline{-2\mathbf{y}^T \mathbf{X}} + \underline{2\theta^T \mathbf{X}^T \mathbf{X}} \right) \\ &= \frac{1}{\sigma^2} \left(-\mathbf{y}^T \mathbf{X} + \theta^T \mathbf{X}^T \mathbf{X} \right) \in \mathbb{R}^{1 \times D}\end{aligned}$$

Maximum Likelihood Estimator

- Setting $\frac{d\mathcal{L}}{d\theta} = \mathbf{0}^T$ results in maximum likelihood parameters θ^*

$$\frac{d\mathcal{L}}{d\theta} = \mathbf{0}^T \iff \cancel{\frac{1}{\sigma^2}}(-\mathbf{y}^T \mathbf{X} + \theta^T \mathbf{X}^T \mathbf{X}) = \mathbf{0}^T$$

$$\iff -\mathbf{y}^T \mathbf{X} + \theta^T \mathbf{X}^T \mathbf{X} = \mathbf{0}^T$$

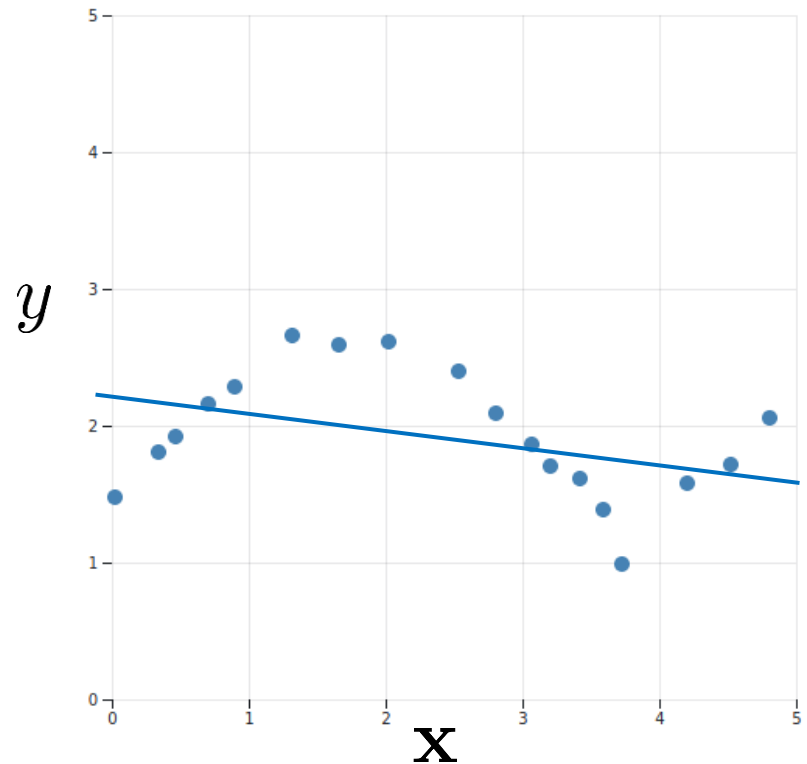
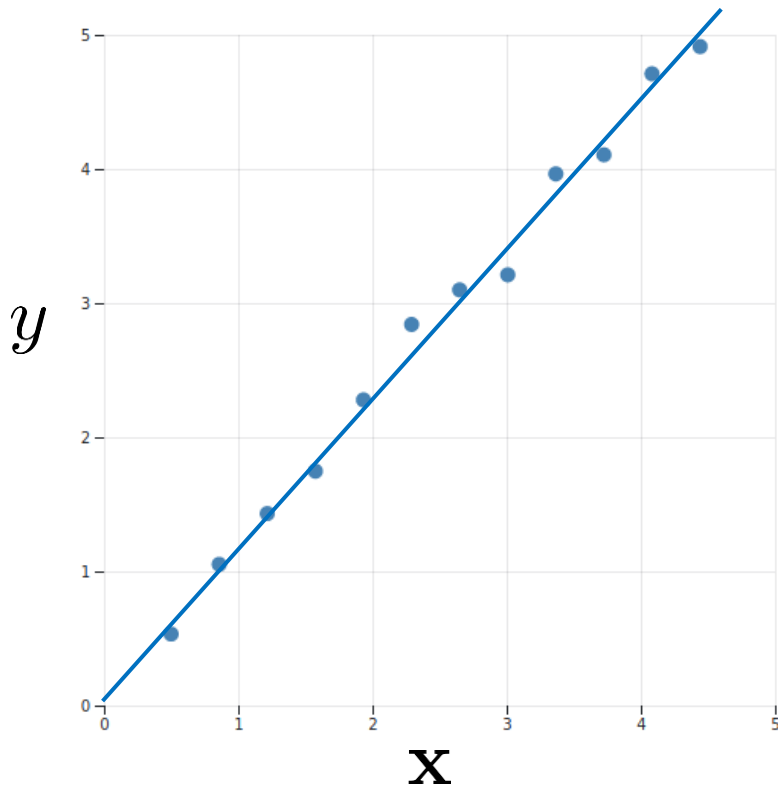
$$\iff \theta^T \mathbf{X}^T \mathbf{X} = \cancel{\mathbf{0}^T} + \mathbf{y}^T \mathbf{X}$$

$$\iff \theta^T = \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \quad (\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

$$\iff \theta = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}$$

Normal Equation

Example: Linear Fit



- Closed form solution for maximum likelihood estimate enables to fit lines

Non-linear Functions

- Linear regression is *linear in parameters*
- We can apply non-linear transformation:

$$f(\mathbf{x}) = \mathbf{x}^T \theta \longrightarrow f(\mathbf{x}) = \phi(\mathbf{x})^T \theta$$

- Let $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^K$ and define $\Phi \in \mathbb{R}^{N \times K}$ as

$$\Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{pmatrix} \in \mathbb{R}^{N \times K}$$

- Everything else stays the same! Use normal equation with Φ instead of \mathbf{X}

$$\theta = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

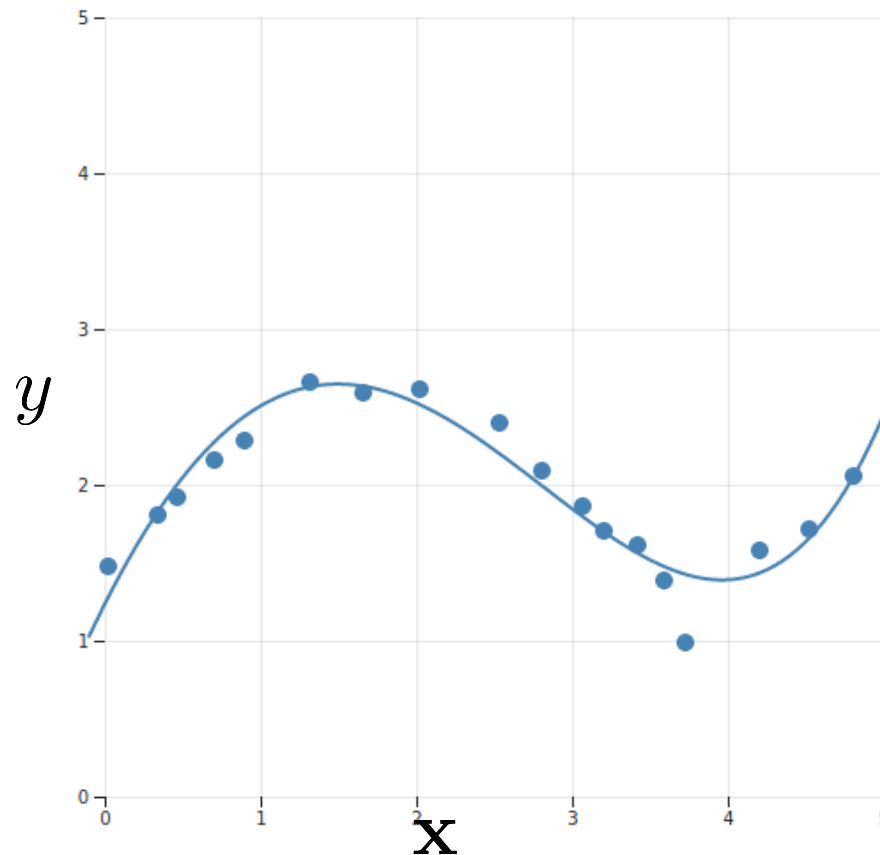
Example: Polynomial transformation

- With polynomial transformation, we can fit polynomials of degree K

$$\phi_{\text{poly}}(\mathbf{x}_n) = \begin{pmatrix} 1 \\ x_1 \\ x_1^2 \\ \vdots \\ x_1^K \\ \vdots \end{pmatrix} \in \mathbb{R}^{DK+1}$$

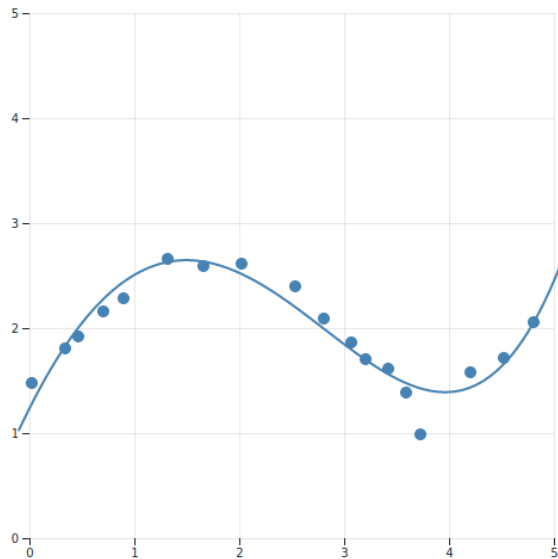
- With K=1 it's “vanilla” linear regression

Example: Polynomial Fit

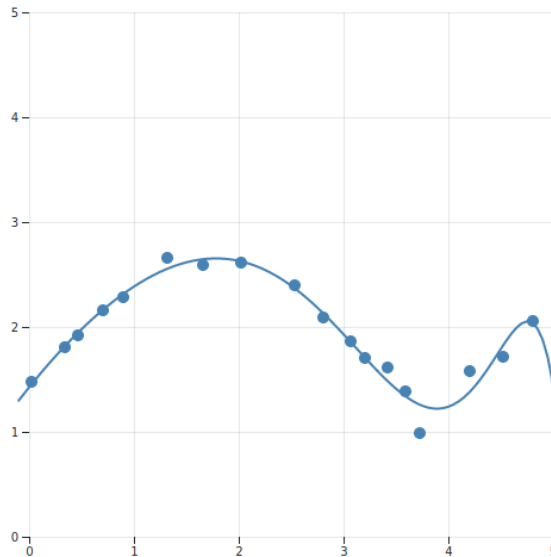


- With a polynomial of degree 3, we get a good fit. But can we do better with higher degrees?

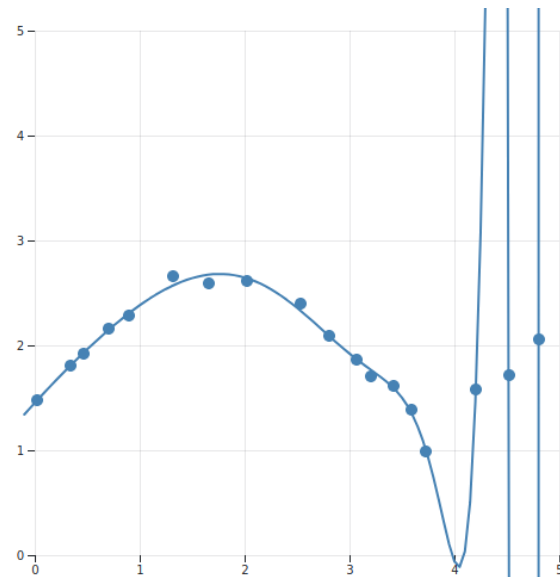
Example: Overfitting



degree 3



degree 5



degree 8

- Increasing the degree, we will get better training error
- But function will have implausible shape

Overfitting and Generalization

- Maximum Likelihood estimates overfit to training data and overconfident predictions
- High capacity models ($K > 5$) tend to fit training data too well
- Usually caused by too large parameter values
- **Solution:** Ensure that parameters do not get too large!

Common Parameter Estimation

- Learning is finding parameters of

$$P(\theta|\mathbf{x}_{1:N}, y_{1:N}) \propto P(y_{1:N}|\mathbf{x}_{1:N}, \theta) \boxed{P(\theta)}$$

- Paradigms for parameter estimation:
 1. Point estimate with uniform prior
→ Maximum Likelihood Estimation
 2. Point estimate with given prior
→ Maximum A posteriori Estimation (MAP)
 3. Determine posterior over the parameters
→ Bayesian Estimation

NLL with Prior

- Assume Gaussian prior for parameters:

$$P(\theta) = \mathcal{N}(\theta|0, b^2\mathbf{I})$$

- NLL is then

$$\mathcal{L}_{\text{MAP}}(\theta) = -\log \prod_{n=1}^N P(y_n|\mathbf{x}_n, \theta) - \log P(\theta)$$

- Inserting Gaussians results in

$$\begin{aligned}\mathcal{L}_{MAP}(\theta) &= \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \phi(\mathbf{x}_n)^T \theta)^2 + \frac{1}{2b^2} \theta^T \theta + \text{const} \\ &= \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{\Phi}\theta)^T (\mathbf{y} - \mathbf{\Phi}\theta) + \frac{1}{2b^2} \theta^T \theta + \text{const}\end{aligned}$$

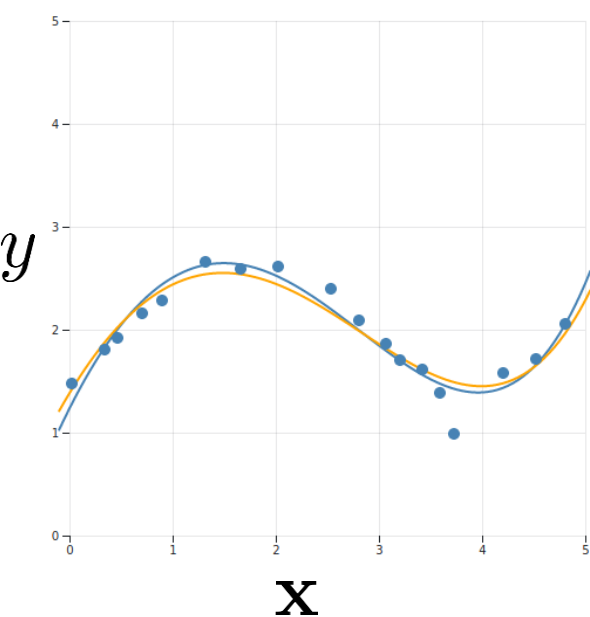
Gradient of MAP NLL

- Same receipt as before:

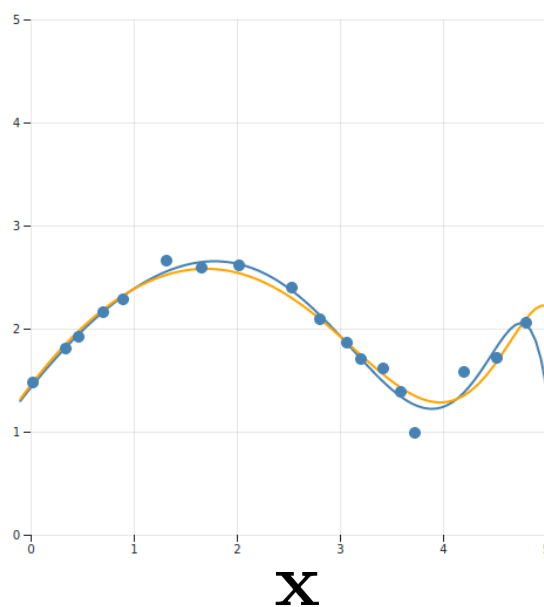
$$\begin{aligned}\frac{d\mathcal{L}_{\text{MAP}}}{d\theta} = \mathbf{0}^T &\iff \frac{1}{\sigma^2} \left(\theta^T \Phi^T \Phi - \mathbf{y}^T \Phi \right) + \frac{1}{b^2} \theta^T = \mathbf{0}^T \\ &\iff \frac{1}{\sigma^2} \theta^T \left(\Phi^T \Phi + \frac{\sigma^2}{b^2} \mathbf{I} \right) - \frac{1}{\sigma^2} \mathbf{y}^T \Phi = \mathbf{0}^T \\ &\iff \theta^T \left(\Phi^T \Phi + \frac{\sigma^2}{b^2} \mathbf{I} \right) = \mathbf{y}^T \Phi \\ &\iff \theta^T = \mathbf{y}^T \Phi \left(\Phi^T \Phi + \frac{\sigma^2}{b^2} \mathbf{I} \right)^{-1} \\ &\iff \theta = \left(\Phi^T \Phi + \frac{\sigma^2}{b^2} \mathbf{I} \right)^{-1} \Phi^T \mathbf{y}\end{aligned}$$

Example: ML vs. MAP Estimate

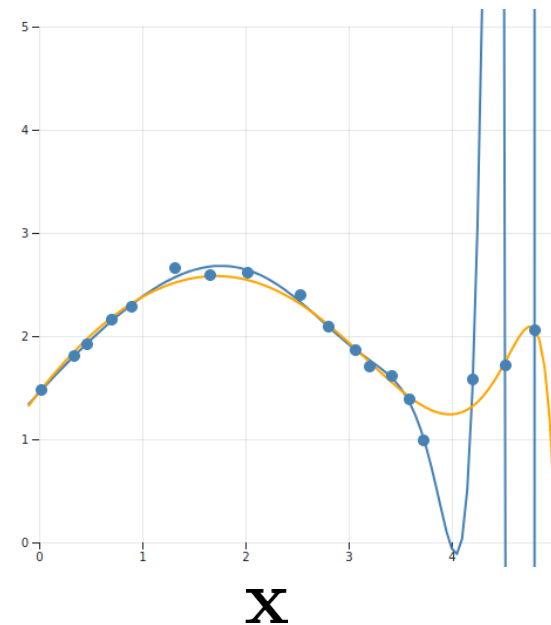
— ML Estimate — MAP Estimate



degree 3



degree 5



degree 8

- Smoother functions even at higher degrees!

Common Parameter Estimation

- Learning is finding parameters of

$$P(\theta|\mathbf{x}_{1:N}, y_{1:N}) \propto P(y_{1:N}|\mathbf{x}_{1:N}, \theta) \quad P(\theta)$$

- Paradigms for parameter estimation:
 1. Point estimate with uniform prior
→ Maximum Likelihood Estimation
 2. Point estimate with given prior
→ Maximum A posteriori Estimation (MAP)
 3. Determine posterior over the parameters
→ Bayesian Estimation

Bayesian Approach

- ML and MAP estimates provide **most likely** parameters after observing the training data

$$\theta^* = \arg \max_{\theta} P(\theta | \mathbf{x}_{1:N}, y_{1:N})$$

- **Bayesian approach**: Posterior over parameters after observing training data

$$\begin{aligned} P(\theta | \mathbf{x}_{1:N}, y_{1:N}) &= \frac{P(y_{1:N} | \mathbf{x}_{1:N}, \theta) P(\theta)}{P(y_{1:N} | \mathbf{x}_{1:N})} \\ &= \frac{P(y_{1:N} | \mathbf{x}_{1:N}, \theta) P(\theta)}{\int P(y_{1:N} | \mathbf{x}_{1:N}, \theta) P(\theta) \, d\theta} \end{aligned}$$

Posterior Predictions

- Predictions $y_* \in \mathbb{R}$ for unseen examples $\mathbf{x}_* \in \mathbb{R}^D$ are weighted average over all possible parameters:

$$P(y_* | x_*, \mathbf{x}_{1:N}, y_{1:N}) = \int P(y_* | \mathbf{x}_*, \theta) P(\theta | \mathbf{x}_{1:N}, y_{1:N}) d\theta$$

- **Advantage:** We take the uncertainty of the parameters into account

Bayesian Linear Regression

- Stick with Gaussians for likelihood and prior

$$P(y_{1:N}|\mathbf{x}_{1:N}, \theta) = \mathcal{N}(\mathbf{y}|\mathbf{\Phi}\theta, \sigma^2\mathbf{I})$$

$$P(\theta) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$$

- To derive the posterior, we can exploit that marginal is product of likelihood and prior

$$P(\theta|y_{1:N}, \mathbf{x}_{1:N}) = \frac{P(y_{1:N}|\mathbf{x}_{1:N}, \theta)P(\theta)}{\int P(y_{1:N}|\mathbf{x}_{1:N}, \theta)P(\theta) d\theta}$$

Inserting Likelihood and Prior

- We insert our choices for likelihood and prior into the numerator:

$$P(y_{1:N}|\mathbf{x}_{1:N}, \theta)P(\theta) = \mathcal{N}(\underline{\mathbf{y}}|\underline{\Phi}\theta, \sigma^2\mathbf{I})\mathcal{N}(\underline{\theta}|\mathbf{m}_0, \mathbf{S}_0)$$

- Nearly a product of two Gaussians, but different variables
- First ensure that Likelihood in terms of θ

Gaussians: Change of variables

- Consider normal distribution of \mathbf{x} with mean as linear function $\mathbf{A}\mathbf{y} + \mathbf{b}$: $\mathcal{N}(\mathbf{x}|\mathbf{A}\mathbf{y} + \mathbf{b}, \Sigma)$
- Change of variable such that normal distribution in terms of \mathbf{y}

$$\mathcal{N}(\mathbf{x}|\mathbf{A}\mathbf{y} + \mathbf{b}, \Sigma) = \eta^{-1} \mathcal{N}(\mathbf{y}|\mathbf{A}'\mathbf{x} + \mathbf{b}', \Sigma')$$

with

$$\Sigma' = (\mathbf{A}^T \Sigma^{-1} \mathbf{A})^{-1}$$

$$\mathbf{A}' = \Sigma' \mathbf{A}^T \Sigma^{-1}$$

$$\mathbf{b}' = -\mathbf{A}'\mathbf{b}$$

$$\eta = |\mathbf{A}|$$

Deriving the Posterior (I)

- Change variables for $\mathcal{N}(\mathbf{y}|\Phi\theta, \sigma^2\mathbf{I})$ gives:

$$\mathcal{N}(\mathbf{y}|\Phi\theta, \sigma^2\mathbf{I}) = \eta_1^{-1} \mathcal{N}(\theta|\mathbf{A}'\mathbf{y}, \Sigma')$$

with

$$\Sigma' = (\sigma^{-2}\Phi^T\Phi)^{-1}$$

$$\mathbf{A}' = (\sigma^{-2}\Phi^T\Phi)^{-1}\sigma^{-2}\Phi^T$$

$$\Sigma' = (\mathbf{A}^T\Sigma^{-1}\mathbf{A})^{-1}$$

$$\mathbf{A}' = \Sigma'\mathbf{A}^T\Sigma^{-1}$$

- Next we want to compute

$$\mathcal{N}(\theta|\mathbf{A}'\mathbf{y}, \Sigma')\mathcal{N}(\theta|\mathbf{m}_0, \mathbf{S}_0)$$

Gaussian Product

- Product of Gaussians is again a (unnormalized) Gaussian given by:

$$\mathcal{N}(\mathbf{x}|\mathbf{a}, \mathbf{A})\mathcal{N}(\mathbf{x}|\mathbf{b}, \mathbf{B}) = \eta^{-1}\mathcal{N}(\mathbf{x}|\mathbf{c}, \mathbf{C}),$$

with

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}$$

$$\mathbf{c} = \mathbf{C} (\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b})$$

$$\eta = \mathcal{N}(\mathbf{a}|\mathbf{b}, \mathbf{A} + \mathbf{B}) = \mathcal{N}(\mathbf{b}|\mathbf{a}, \mathbf{A} + \mathbf{B})$$

Deriving the Posterior (II)

- Using the product we can simplify further

$$\mathcal{N}(\mathbf{y}|\Phi\theta, \sigma^2\mathbf{I}) \mathcal{N}(\theta|\mathbf{m}_0, \mathbf{S}_0)$$

$$= \eta_1^{-1} \mathcal{N}(\theta | (\sigma^{-2}\Phi^T\Phi)^{-1}\sigma^{-2}\Phi^T\mathbf{y}, (\sigma^{-2}\Phi^T\Phi)^{-1}) \mathcal{N}(\theta|\mathbf{m}_0, \mathbf{S}_0)$$

change of variables

$$= \eta_1^{-1}\eta_2^{-1} \mathcal{N}(\theta|\mathbf{m}_N, \mathbf{S}_N)$$

with

$$\begin{aligned}\mathbf{C} &= (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} \\ \mathbf{c} &= \mathbf{C} (\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b})\end{aligned}$$

$$\mathbf{S}_N = (\sigma^{-2}\Phi^T\Phi + \mathbf{S}_0^{-1})^{-1}$$

$$\begin{aligned}\mathbf{m}_N &= \mathbf{S}_N (\sigma^{-2}\Phi^T\Phi (\sigma^{-2}\Phi^T\Phi)^{-1}\sigma^{-2}\Phi^T\mathbf{y} + \mathbf{S}_0^{-1}\mathbf{m}_0) \\ &= \mathbf{S}_N (\sigma^{-2}\Phi^T\mathbf{y} + \mathbf{S}_0^{-1}\mathbf{m}_0)\end{aligned}$$

Solving the Marginal Likelihood

- The numerator of the posterior is

$$P(y_{1:N}|\theta, \mathbf{x}_{1:N})P(\theta) = \eta_1^{-1}\eta_2^{-1}\mathcal{N}(\theta|\mathbf{m}_N, \mathbf{S}_N)$$

- The denominator is

$$\begin{aligned} P(y_{1:N}|\mathbf{x}_{1:N}) &= \int P(y_{1:N}|\theta, \mathbf{x}_{1:N})P(\theta) \, d\theta \\ &= \int \eta_1^{-1}\eta_2^{-1}\mathcal{N}(\theta|\mathbf{m}_N, \mathbf{S}_N) \, d\theta \\ &= \eta_1^{-1}\eta_2^{-1} \underbrace{\int \mathcal{N}(\theta|\mathbf{m}_N, \mathbf{S}_N) \, d\theta}_{=1} \\ &= \eta_1^{-1}\eta_2^{-1} \end{aligned}$$

Posterior is Gaussian

- Putting everything together, gives us:

$$P(\theta|y_{1:N}, \mathbf{x}_{1:N}) = \mathcal{N}(\theta|\mathbf{m}_N, \mathbf{S}_N)$$

with

$$\begin{aligned}\mathbf{S}_N &= (\sigma^{-2} \mathbf{\Phi}^T \mathbf{\Phi} + \mathbf{S}_0^{-1})^{-1} \\ \mathbf{m}_N &= \mathbf{S}_N (\sigma^{-2} \mathbf{\Phi}^T \mathbf{y} + \mathbf{S}_0^{-1} \mathbf{m}_0)\end{aligned}$$

- Posterior has the same form as the prior.

Conjugate Prior

- Prior is a **conjugate Prior** for a likelihood function if the posterior is of the *same* form/type as the prior.
- Gaussians with known Σ are self-conjugate.
- Conjugate priors lead to closed-form solutions
- Posterior is prior with updated parameters

Posterior Prediction

- Using the posterior, we can finally make predictions:

$$P(y_*|\mathbf{x}_{1:N}, y_{1:N}) = \int P(y_*|\mathbf{x}_*, \theta) P(\theta|\mathbf{x}_{1:N}, y_{1:N}) \, d\theta$$

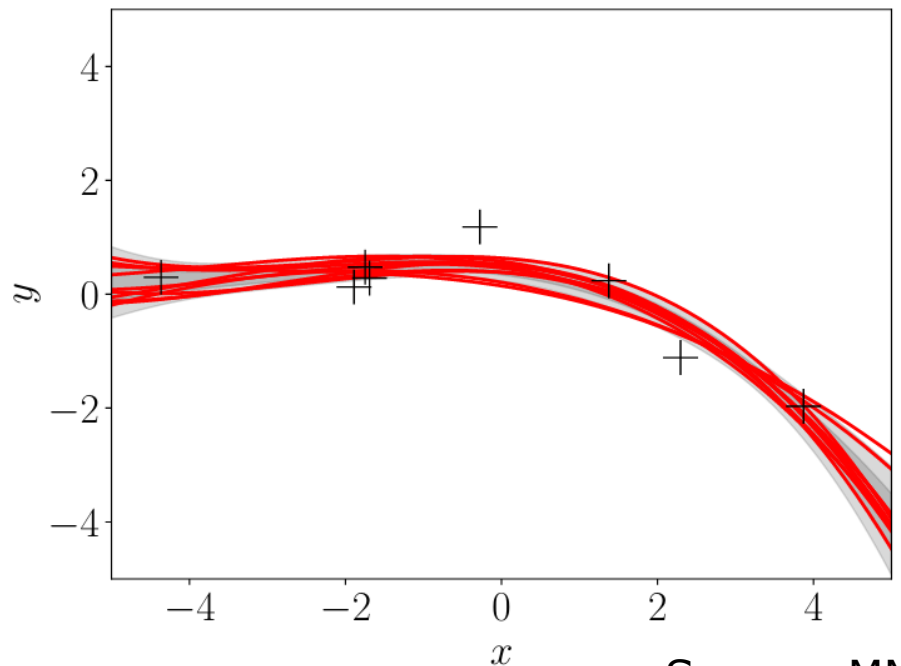
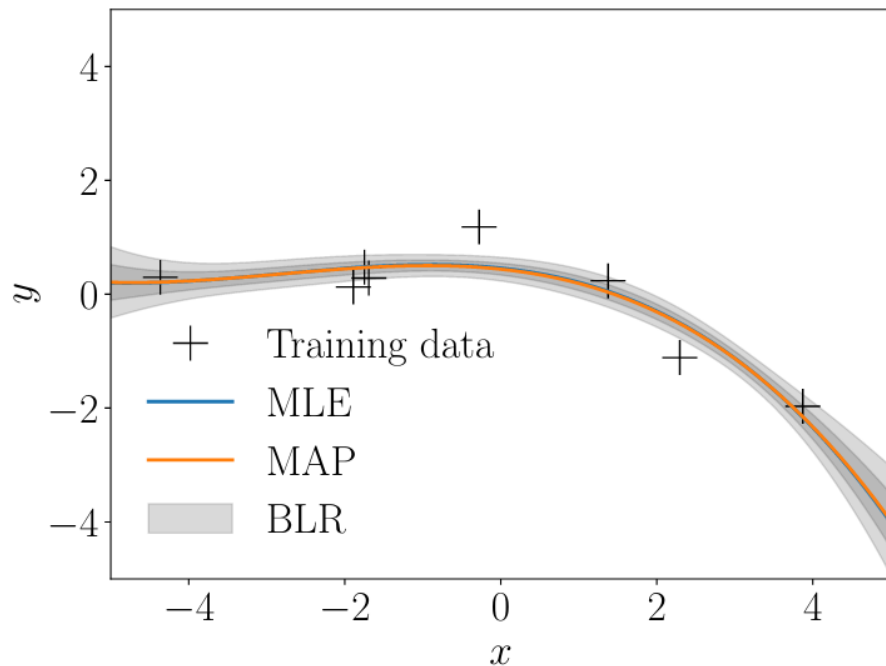
- Inserting likelihood and posterior gives:

$$= \int \mathcal{N}(y_*|\phi(\mathbf{x}_*)\theta, \sigma^2) \mathcal{N}(\theta|\mathbf{m}_N, \mathbf{S}_N) \, d\theta$$

- Product of Gaussians is Gaussian again. One can derive that

$$= \mathcal{N}(y_*|\phi^T(\mathbf{x}_*)\mathbf{m}_N, \phi^T(\mathbf{x}_*)\mathbf{S}_N\phi(\mathbf{x}_*) + \sigma^2)$$

Example: Bayesian Regression

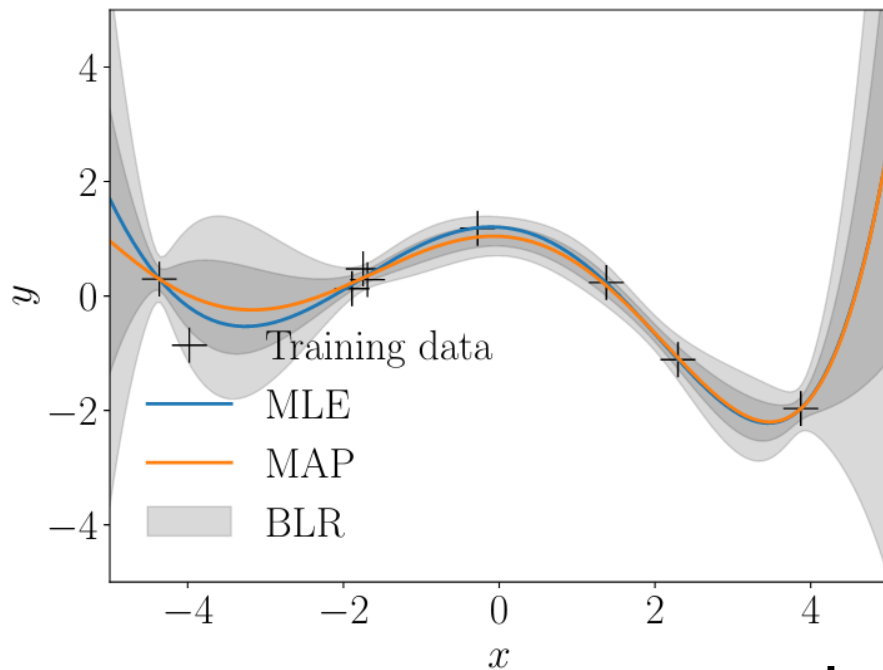


Source: MML

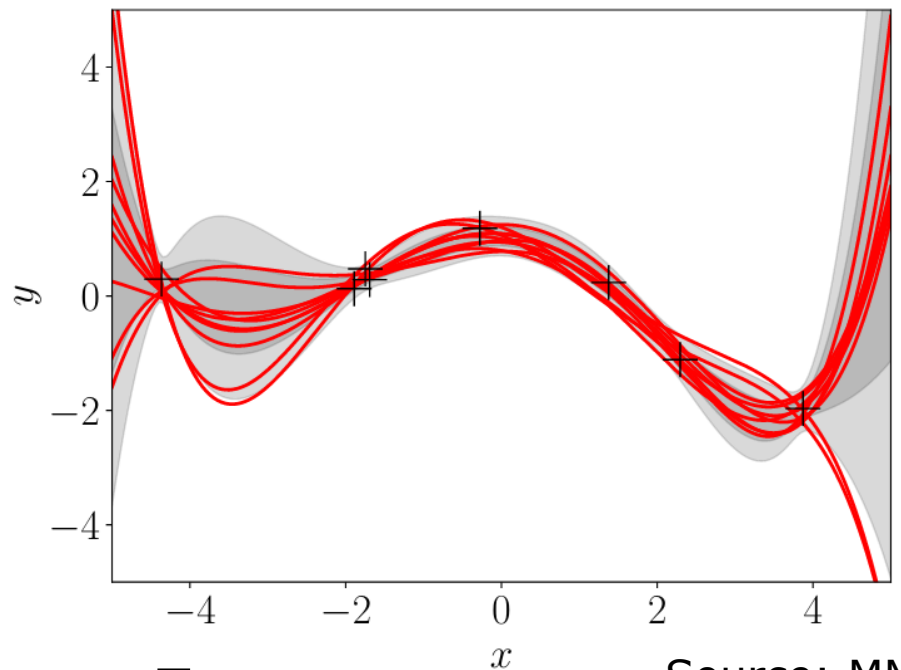
degree 3

- With higher degrees, we see also larger confidence intervals.

Example: Bayesian Regression



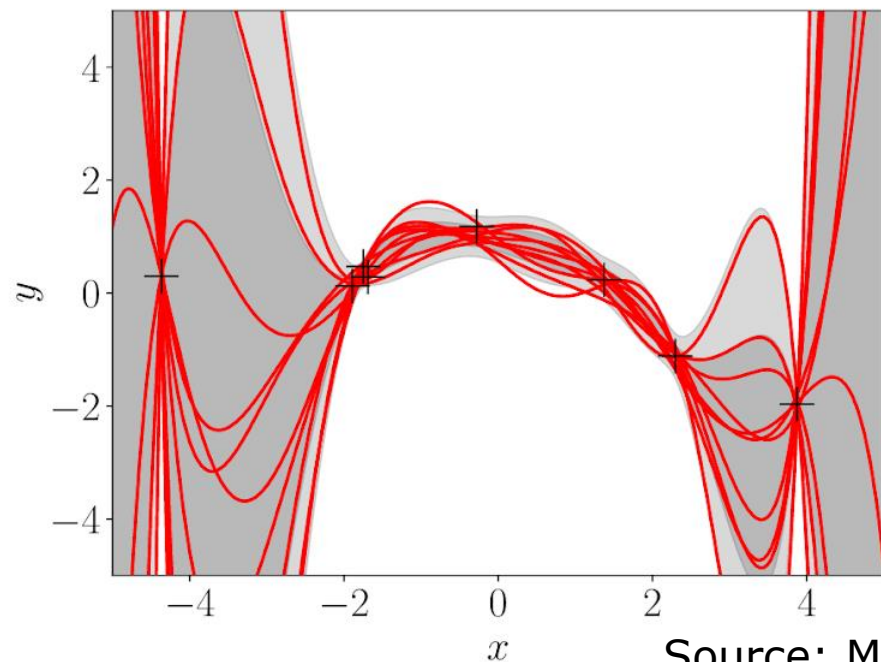
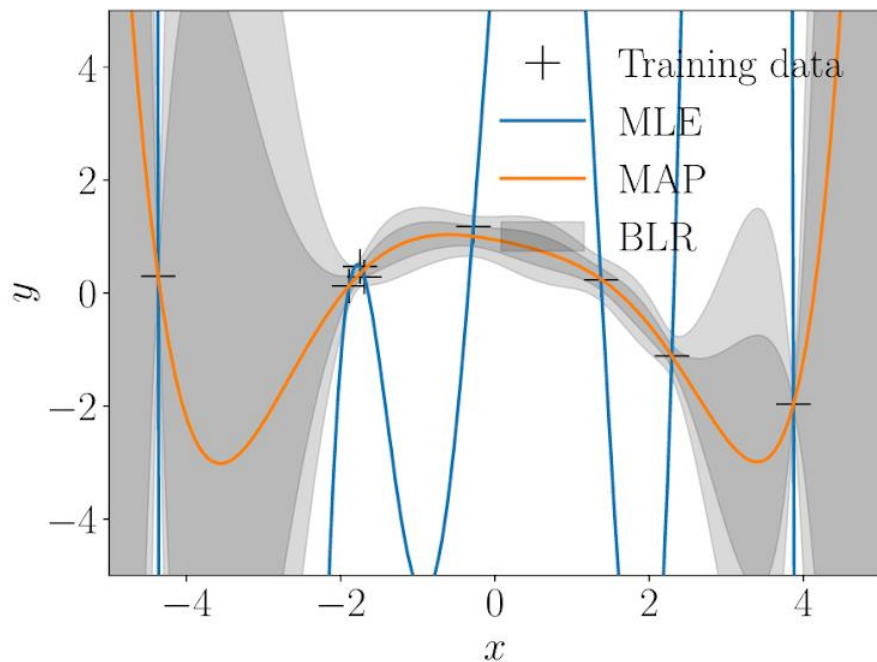
degree 5



Source: MML

- With higher degrees, we see also larger confidence intervals.

Example: Bayesian Regression



degree 7

- With higher degrees, we see also larger confidence intervals.

Why Bayesian Regression?

- In robotics several properties that are advantageous:
 - 1. Uncertainty** with variance of predictions
 - 2.** Determine "regions of uncertainty"
 - 3. Incremental learning** possible: posterior can be prior for next round of learning!

Discriminative and Generative Models

- With Bayes Theorem, we can express as follows:

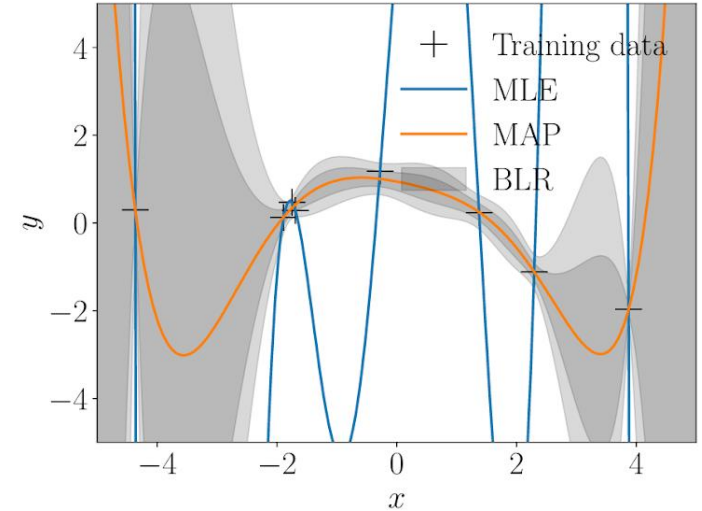
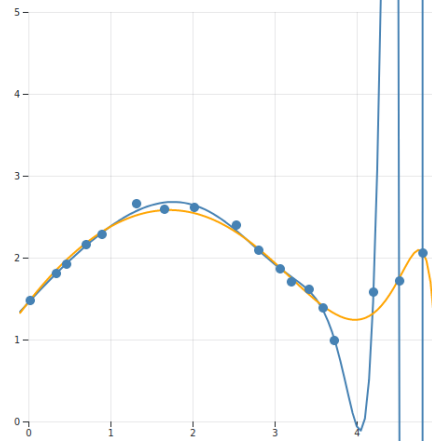
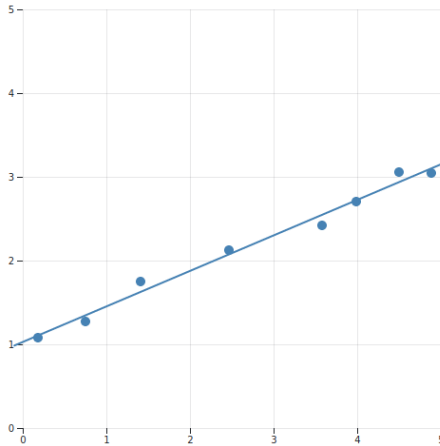
$$P(y|\mathbf{x}) = \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} = \frac{P(\mathbf{x}|y)P(y)}{\int P(\mathbf{x}|y)P(y) \, dy}$$

- Up to now, we used a **discriminative model** we modeled directly $P(y|\mathbf{x})$
- But we could equally model $P(\mathbf{x}, y)$ or $P(\mathbf{x}|y), P(y)$ to get $P(y|\mathbf{x}) \rightarrow$ **generative model**

Which model to use?

- No definite answer (depends)
- Points to consider:
 1. Inference simpler with discriminative models
 2. The data x is of higher dimension than the label; needs complex models
 3. $P(x|y)$ allows to directly model the generation process and generate data.
 4. Generative models can handle missing data.
 5. Priors can incorporate expert knowledge in the prior in generative models.

Summary



- Linear Regression and variants
- Bayesian Linear Regression
- Discriminative vs. Generative Models

See you next week!