

Probabilistic Models

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

Machine Learning

Outline

- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation
- 4 Bayesian Estimation and Inference*
 - Gaussian Process

Outline

- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation
- 4 Bayesian Estimation and Inference*
 - Gaussian Process

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$
- Model \mathbb{F} : a collection of functions parametrized by Θ

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$
- Model \mathbb{F} : a collection of functions parametrized by Θ
- Goal: to train a function f such that, given a new data point \mathbf{x}' , the output value

$$\hat{\mathbf{y}} = f(\mathbf{x}'; \Theta)$$

is closest to the correct label \mathbf{y}'

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$
- Model \mathbb{F} : a collection of functions parametrized by Θ
- Goal: to train a function f such that, given a new data point \mathbf{x}' , the output value

$$\hat{\mathbf{y}} = f(\mathbf{x}'; \Theta)$$

is closest to the correct label \mathbf{y}'

- Examples in \mathbb{X} are usually assumed to be i.i.d. sampled from random variables (\mathbf{x}, \mathbf{y}) following some data generating distribution $P(\mathbf{x}, \mathbf{y})$

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$
- Model \mathbb{F} : a collection of functions parametrized by Θ
- Goal: to train a function f such that, given a new data point \mathbf{x}' , the output value

$$\hat{\mathbf{y}} = f(\mathbf{x}'; \Theta)$$

is closest to the correct label \mathbf{y}'

- Examples in \mathbb{X} are usually assumed to be i.i.d. sampled from random variables (\mathbf{x}, \mathbf{y}) following some data generating distribution $P(\mathbf{x}, \mathbf{y})$
- In probabilistic models, we write $f(\mathbf{x}'; \Theta)$ as $P(\mathbf{y} = \mathbf{y} | \mathbf{x} = \mathbf{x}')$ and a prediction is made by:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} P(\mathbf{y} = \mathbf{y} | \mathbf{x} = \mathbf{x}'; \Theta)$$

Predictions based on Probability

- Supervised learning, we are given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$
- Model \mathbb{F} : a collection of functions parametrized by Θ
- Goal: to train a function f such that, given a new data point \mathbf{x}' , the output value

$$\hat{\mathbf{y}} = f(\mathbf{x}'; \Theta)$$

is closest to the correct label \mathbf{y}'

- Examples in \mathbb{X} are usually assumed to be i.i.d. sampled from random variables (\mathbf{x}, \mathbf{y}) following some data generating distribution $P(\mathbf{x}, \mathbf{y})$
- In probabilistic models, we write $f(\mathbf{x}'; \Theta)$ as $P(\mathbf{y} = \mathbf{y} | \mathbf{x} = \mathbf{x}')$ and a prediction is made by:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} P(\mathbf{y} = \mathbf{y} | \mathbf{x} = \mathbf{x}'; \Theta)$$

- How to find Θ ?

Function (Θ) as Point Estimate

- Regard $\Theta(f)$ as an estimate of the “true” $\Theta^*(f^*)$
 - Mapped from the training set \mathbb{X}

Function (Θ) as Point Estimate

- Regard $\Theta(f)$ as an estimate of the “true” $\Theta^*(f^*)$
 - Mapped from the training set \mathbb{X}
- *Maximum a posteriori (MAP) estimation:*

$$\Theta_{\text{MAP}} = \arg \max_{\Theta} P(\Theta | \mathbb{X}) = \arg \max_{\Theta} \frac{P(\mathbb{X} | \Theta)P(\Theta)}{P(\mathbb{X})} = \arg \max_{\Theta} P(\mathbb{X} | \Theta)P(\Theta)$$

- Θ is regarded as a random variable

Function (Θ) as Point Estimate

- Regard $\Theta(f)$ as an estimate of the “true” $\Theta^*(f^*)$
 - Mapped from the training set \mathbb{X}
- *Maximum a posteriori (MAP) estimation:*

$$\Theta_{\text{MAP}} = \arg \max_{\Theta} P(\Theta | \mathbb{X}) = \arg \max_{\Theta} \frac{P(\mathbb{X} | \Theta)P(\Theta)}{P(\mathbb{X})} = \arg \max_{\Theta} P(\mathbb{X} | \Theta)P(\Theta)$$

- Θ is regarded as a random variable
- *Maximum likelihood (ML) estimation:*

$$\Theta_{\text{ML}} = \arg \max_{\Theta} P(\mathbb{X} | \Theta)$$

- Assumes a uniform $P(\Theta)$ (i.e., does not prefer a particular Θ)

Function (Θ) as Point Estimate

- Regard $\Theta(f)$ as an estimate of the “true” $\Theta^*(f^*)$
 - Mapped from the training set \mathbb{X}
- *Maximum a posteriori (MAP) estimation:*

$$\Theta_{\text{MAP}} = \arg \max_{\Theta} P(\Theta | \mathbb{X}) = \arg \max_{\Theta} \frac{P(\mathbb{X} | \Theta)P(\Theta)}{P(\mathbb{X})} = \arg \max_{\Theta} P(\mathbb{X} | \Theta)P(\Theta)$$

- Θ is regarded as a random variable
- *Maximum likelihood (ML) estimation:*

$$\Theta_{\text{ML}} = \arg \max_{\Theta} P(\mathbb{X} | \Theta)$$

- Assumes a uniform $P(\Theta)$ (i.e., does not prefer a particular Θ)
- After being solved, $\Theta_{\text{ML}/\text{MAP}}$ is treated as a constant when make a prediction $\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta_{\text{ML}/\text{MAP}})$

Outline

1 Probabilistic Models

2 Maximum Likelihood Estimation

- Linear Regression
- Logistic Regression

3 Maximum A Posteriori Estimation

4 Bayesian Estimation and Inference*

- Gaussian Process

Outline

1 Probabilistic Models

2 Maximum Likelihood Estimation

- Linear Regression
- Logistic Regression

3 Maximum A Posteriori Estimation

4 Bayesian Estimation and Inference*

- Gaussian Process

Probability Interpretation

- Assumption: $y = f^*(\mathbf{x}) + \varepsilon$, where
 - $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$

Probability Interpretation

- Assumption: $y = f^*(\mathbf{x}) + \varepsilon$, where
 - $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$
 - $f^*(\mathbf{x}; \mathbf{w}^*) = \mathbf{w}^{*\top} \mathbf{x}$ is a deterministic function
 - All variables in \mathbf{x} are z -normalized, so there's no bias term b in f^*

Probability Interpretation

- Assumption: $y = f^*(\mathbf{x}) + \varepsilon$, where
 - $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$
 - $f^*(\mathbf{x}; \mathbf{w}^*) = \mathbf{w}^{*\top} \mathbf{x}$ is a deterministic function
 - All variables in \mathbf{x} are z -normalized, so there's no bias term b in f^*
- We have $(y | \mathbf{x} = \mathbf{x}) \sim \mathcal{N}(\mathbf{w}^{*\top} \mathbf{x}, \beta^{-1})$

Probability Interpretation

- Assumption: $y = f^*(\mathbf{x}) + \varepsilon$, where
 - $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$
 - $f^*(\mathbf{x}; \mathbf{w}^*) = \mathbf{w}^{*\top} \mathbf{x}$ is a deterministic function
 - All variables in \mathbf{x} are z -normalized, so there's no bias term b in f^*
- We have $(y | \mathbf{x} = \mathbf{x}) \sim \mathcal{N}(\mathbf{w}^{*\top} \mathbf{x}, \beta^{-1})$
- So, our goal is to find \mathbf{w} as close to \mathbf{w}^* as possible such that:

$$\hat{y} = \arg \max_y P(y | \mathbf{x} = \mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \mathbf{x}$$

- Note that \hat{y} is irrelevant to β , so we don't need to solve β

Probability Interpretation

- Assumption: $y = f^*(\mathbf{x}) + \varepsilon$, where
 - $\varepsilon \sim \mathcal{N}(0, \beta^{-1})$
 - $f^*(\mathbf{x}; \mathbf{w}^*) = \mathbf{w}^{*\top} \mathbf{x}$ is a deterministic function
 - All variables in \mathbf{x} are z -normalized, so there's no bias term b in f^*
- We have $(y | \mathbf{x} = \mathbf{x}) \sim \mathcal{N}(\mathbf{w}^{*\top} \mathbf{x}, \beta^{-1})$
- So, our goal is to find \mathbf{w} as close to \mathbf{w}^* as possible such that:

$$\hat{y} = \arg \max_y P(y | \mathbf{x} = \mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \mathbf{x}$$

- Note that \hat{y} is irrelevant to β , so we don't need to solve β
- ML estimation for \mathbf{w}^* :

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

ML Estimation I

- Problem:

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$P(\mathbb{X} | \mathbf{w}) = \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w})$$

ML Estimation I

- Problem:

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned} P(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) P(\mathbf{x}^{(i)}) \end{aligned}$$

ML Estimation I

- Problem:

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} \mathcal{P}(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned} \mathcal{P}(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N \mathcal{P}(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N \mathcal{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathcal{P}(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N \mathcal{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathcal{P}(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) \mathcal{P}(\mathbf{x}^{(i)}) \\ &= \prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) \mathcal{P}(\mathbf{x}^{(i)}) \end{aligned}$$

ML Estimation I

- Problem:

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} \mathcal{P}(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned}\mathcal{P}(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N \mathcal{P}(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N \mathcal{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathcal{P}(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N \mathcal{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathcal{P}(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) \mathcal{P}(\mathbf{x}^{(i)}) \\ &= \prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) \mathcal{P}(\mathbf{x}^{(i)})\end{aligned}$$

- To make the problem tractable, we prefer “sums” over “products”

ML Estimation I

- Problem:

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned} P(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) P(\mathbf{x}^{(i)}) \\ &= \prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) P(\mathbf{x}^{(i)}) \end{aligned}$$

- To make the problem tractable, we prefer “sums” over “products”
- We can instead maximize the *log likelihood*

$$\arg \max_{\mathbf{w}} \text{log} P(\mathbb{X} | \mathbf{w})$$

- The optimal point does not change since log is monotone increasing

ML Estimation I

- Problem:

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned} P(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) P(\mathbf{x}^{(i)}) \\ &= \prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) P(\mathbf{x}^{(i)}) \end{aligned}$$

- To make the problem tractable, we prefer “sums” over “products”
- We can instead maximize the **log likelihood**

$$\begin{aligned} \arg \max_{\mathbf{w}} \log P(\mathbb{X} | \mathbf{w}) \\ = \arg \max_{\mathbf{w}} \log \left[\prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) P(\mathbf{x}^{(i)}) \right] \end{aligned}$$

- The optimal point does not change since log is monotone increasing

ML Estimation I

- Problem:

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} \mathbb{P}(\mathbb{X} | \mathbf{w})$$

- Since we assume i.i.d. samples, we have

$$\begin{aligned} \mathbb{P}(\mathbb{X} | \mathbf{w}) &= \prod_{i=1}^N \mathbb{P}(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) = \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)} | \mathbf{w}) \\ &= \prod_{i=1}^N \mathbb{P}(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \mathbb{P}(\mathbf{x}^{(i)}) = \prod_i \mathcal{N}(y^{(i)}; \mathbf{w}^\top \mathbf{x}^{(i)}, \sigma^2) \mathbb{P}(\mathbf{x}^{(i)}) \\ &= \prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) \mathbb{P}(\mathbf{x}^{(i)}) \end{aligned}$$

- To make the problem tractable, we prefer “sums” over “products”
- We can instead maximize the **log likelihood**

$$\begin{aligned} \arg \max_{\mathbf{w}} \log \mathbb{P}(\mathbb{X} | \mathbf{w}) \\ &= \arg \max_{\mathbf{w}} \log \left[\prod_i \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2\right) \mathbb{P}(\mathbf{x}^{(i)}) \right] \\ &= \arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i \log \mathbb{P}(\mathbf{x}^{(i)}) \end{aligned}$$

- The optimal point does not change since log is monotone increasing

ML Estimation II

$$\arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i P(\mathbf{x}^{(i)})$$

- Ignoring terms irrelevant to \mathbf{w} , we have

$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2$$

ML Estimation II

$$\arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i P(\mathbf{x}^{(i)})$$

- Ignoring terms irrelevant to \mathbf{w} , we have

$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2$$

- Effectively, we seek for \mathbf{w} by *minimizing the SSE* (sum of square errors), as we have done before
 - Can solved analytically
 - Or numerically by, e.g., the stochastic gradient descent algorithm

ML Estimation II

$$\arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i P(\mathbf{x}^{(i)})$$

- Ignoring terms irrelevant to \mathbf{w} , we have

$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2$$

- Effectively, we seek for \mathbf{w} by *minimizing the SSE* (sum of square errors), as we have done before
 - Can solved analytically
 - Or numerically by, e.g., the stochastic gradient descent algorithm
- This new perspective explains our ad hoc choice of SSE for empirical risk minimization
 - Checking assumptions helps understand when model works the best

ML Estimation II

$$\arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i P(\mathbf{x}^{(i)})$$

- Ignoring terms irrelevant to \mathbf{w} , we have

$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2$$

- Effectively, we seek for \mathbf{w} by **minimizing the SSE** (sum of square errors), as we have done before
 - Can solved analytically
 - Or numerically by, e.g., the stochastic gradient descent algorithm
- This new perspective explains our ad hoc choice of SSE for empirical risk minimization
 - Checking assumptions helps understand when model works the best
- Also motivates new models

ML Estimation II

$$\arg \max_{\mathbf{w}} N \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 + \sum_i P(\mathbf{x}^{(i)})$$

- Ignoring terms irrelevant to \mathbf{w} , we have

$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} \sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2$$

- Effectively, we seek for \mathbf{w} by **minimizing the SSE** (sum of square errors), as we have done before
 - Can solved analytically
 - Or numerically by, e.g., the stochastic gradient descent algorithm
- This new perspective explains our ad hoc choice of SSE for empirical risk minimization
 - Checking assumptions helps understand when model works the best
- Also motivates new models. Probabilistic model for classification?

Outline

1 Probabilistic Models

2 Maximum Likelihood Estimation

- Linear Regression
- Logistic Regression

3 Maximum A Posteriori Estimation

4 Bayesian Estimation and Inference*

- Gaussian Process

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

- In regression, we assume $(y | \mathbf{x}) \sim \mathcal{N}$ (based on $y = f^*(\mathbf{x}) + \epsilon$)

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

- In regression, we assume $(y | \mathbf{x}) \sim \mathcal{N}$ (based on $y = f^*(\mathbf{x}) + \varepsilon$)
- However, Gaussian distribution is **not** applicable to binary classification
 - The values of y should concentrate in either 1 or -1

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

- In regression, we assume $(y | \mathbf{x}) \sim \mathcal{N}$ (based on $y = f^*(\mathbf{x}) + \varepsilon$)
- However, Gaussian distribution is **not** applicable to binary classification
 - The values of y should concentrate in either 1 or -1
- Which distribution to assume?

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

- In regression, we assume $(y | \mathbf{x}) \sim \mathcal{N}$ (based on $y = f^*(\mathbf{x}) + \varepsilon$)
- However, Gaussian distribution is **not** applicable to binary classification
 - The values of y should concentrate in either 1 or -1
- Which distribution to assume?
- Coin flipping: $(y | \mathbf{x}) \sim \text{Bernoulli}(\rho)$, where

$$P(y | \mathbf{x}; \rho) = \rho^{y'} (1 - \rho)^{(1-y')}, \text{ where } y' = \frac{y+1}{2} \in \{0, 1\}$$

Probabilistic Models for Binary Classification

- Probabilistic models:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \Theta)$$

- In regression, we assume $(y | \mathbf{x}) \sim \mathcal{N}$ (based on $y = f^*(\mathbf{x}) + \varepsilon$)
- However, Gaussian distribution is **not** applicable to binary classification
 - The values of y should concentrate in either 1 or -1
- Which distribution to assume?
- Coin flipping: $(y | \mathbf{x}) \sim \text{Bernoulli}(\rho)$, where

$$P(y | \mathbf{x}; \rho) = \rho^{y'} (1 - \rho)^{(1-y')}, \text{ where } y' = \frac{y+1}{2} \in \{0, 1\}$$

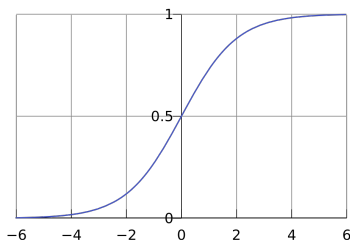
- ML estimate $P(\mathbb{X} | \rho)$? How to relate \mathbf{x} to ρ ?

Logistic Function

- Recall that the *logistic function*

$$\sigma(z) = \frac{\exp(z)}{\exp(z) + 1} = \frac{1}{1 + \exp(-z)}$$

is commonly used as a parametrizing function of the Bernoulli distribution



Logistic Function

- Recall that the *logistic function*

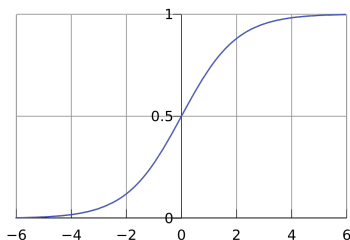
$$\sigma(z) = \frac{\exp(z)}{\exp(z) + 1} = \frac{1}{1 + \exp(-z)}$$

is commonly used as a parametrizing function of the Bernoulli distribution

- We have

$$P(y|\mathbf{x};z) = \sigma(z)^{y'} (1 - \sigma(z))^{(1-y')}$$

- The larger z , the higher chance we get a “positive flip”



Logistic Function

- Recall that the *logistic function*

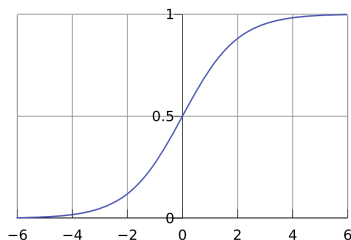
$$\sigma(z) = \frac{\exp(z)}{\exp(z) + 1} = \frac{1}{1 + \exp(-z)}$$

is commonly used as a parametrizing function of the Bernoulli distribution

- We have

$$P(y|\mathbf{x};z) = \sigma(z)^{y'} (1 - \sigma(z))^{(1-y')}$$

- The larger z , the higher chance we get a “positive flip”
- How to relate \mathbf{x} to z ?



Logistic Regression

- In *logistic regression*, we let

$$z = \mathbf{w}^\top \mathbf{x}$$

- Basically, z is the projection of \mathbf{x} along the direction \mathbf{w}

Logistic Regression

- In *logistic regression*, we let

$$z = \mathbf{w}^\top \mathbf{x}$$

- Basically, z is the projection of \mathbf{x} along the direction \mathbf{w}
- We have

$$P(y|\mathbf{x};\mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x})^{y'} [1 - \sigma(\mathbf{w}^\top \mathbf{x})]^{(1-y')}$$

- Prediction:

$$\hat{y} = \arg \max_y P(y|\mathbf{x};\mathbf{w})$$

Logistic Regression

- In *logistic regression*, we let

$$z = \mathbf{w}^\top \mathbf{x}$$

- Basically, z is the projection of \mathbf{x} along the direction \mathbf{w}
- We have

$$P(y|\mathbf{x};\mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x})^{y'} [1 - \sigma(\mathbf{w}^\top \mathbf{x})]^{(1-y')}$$

- Prediction:

$$\hat{y} = \arg \max_y P(y|\mathbf{x};\mathbf{w}) = \text{sign}(\mathbf{w}^\top \mathbf{x})$$

Logistic Regression

- In *logistic regression*, we let

$$z = \mathbf{w}^\top \mathbf{x}$$

- Basically, z is the projection of \mathbf{x} along the direction \mathbf{w}
- We have

$$P(y|\mathbf{x};\mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x})^{y'} [1 - \sigma(\mathbf{w}^\top \mathbf{x})]^{(1-y')}$$

- Prediction:

$$\hat{y} = \arg \max_y P(y|\mathbf{x};\mathbf{w}) = \text{sign}(\mathbf{w}^\top \mathbf{x})$$

- How to learn \mathbf{w} from \mathbb{X} ?

Logistic Regression

- In *logistic regression*, we let

$$z = \mathbf{w}^\top \mathbf{x}$$

- Basically, z is the projection of \mathbf{x} along the direction \mathbf{w}
- We have

$$P(y|\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x})^{y'} [1 - \sigma(\mathbf{w}^\top \mathbf{x})]^{(1-y')}$$

- Prediction:

$$\hat{y} = \arg \max_y P(y|\mathbf{x}; \mathbf{w}) = \text{sign}(\mathbf{w}^\top \mathbf{x})$$

- How to learn \mathbf{w} from \mathbb{X} ?
- ML estimation:

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

ML Estimation

- Log-likelihood:

$$\begin{aligned}\log P(\mathbb{X} | \mathbf{w}) &= \log \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) \\ &= \log \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w})\end{aligned}$$

ML Estimation

- Log-likelihood:

$$\begin{aligned}\log P(\mathbb{X} | \mathbf{w}) &= \log \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) \\ &= \log \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w}) \\ &\propto \log \prod_i \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})^{y^{(i)}} [1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})]^{(1-y^{(i)})}\end{aligned}$$

ML Estimation

- Log-likelihood:

$$\begin{aligned}\log P(\mathbb{X} | \mathbf{w}) &= \log \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) \\ &= \log \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w}) \\ &\propto \log \prod_i \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})^{y^{(i)}} [1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})]^{(1-y^{(i)})} \\ &= \sum_i y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)} - \log(1 + e^{\mathbf{w}^\top \mathbf{x}^{(i)}}) \text{ [Homework]}\end{aligned}$$

- Unlike in linear regression, we cannot solve \mathbf{w} analytically in a closed form via

$$\nabla_{\mathbf{w}} \log P(\mathbb{X} | \mathbf{w}) = \sum_{t=1}^N [y^{(i)} - \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})] \mathbf{x}^{(i)} = \mathbf{0}$$

ML Estimation

- Log-likelihood:

$$\begin{aligned}\log P(\mathbb{X} | \mathbf{w}) &= \log \prod_{i=1}^N P(\mathbf{x}^{(i)}, y^{(i)} | \mathbf{w}) \\ &= \log \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{x}^{(i)} | \mathbf{w}) \\ &\propto \log \prod_i \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})^{y^{(i)}} [1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})]^{(1-y^{(i)})} \\ &= \sum_i y^{(i)} \mathbf{w}^\top \mathbf{x}^{(i)} - \log(1 + e^{\mathbf{w}^\top \mathbf{x}^{(i)}}) \text{ [Homework]}\end{aligned}$$

- Unlike in linear regression, we cannot solve \mathbf{w} analytically in a closed form via

$$\nabla_{\mathbf{w}} \log P(\mathbb{X} | \mathbf{w}) = \sum_{i=1}^N [y^{(i)} - \sigma(\mathbf{w}^\top \mathbf{x}^{(i)})] \mathbf{x}^{(i)} = \mathbf{0}$$

- But since $\log P(\mathbb{X} | \mathbf{w})$ is differentiable w.r.t. \mathbf{w} , we can solve \mathbf{w}_{ML}^* numerically using stochastic gradient descent (SGD)
 - It can be shown that $\log P(\mathbb{X} | \mathbf{w})$ is concave in terms of \mathbf{w} [1]
 - SGD finds global optimal

Outline

- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation**
- 4 Bayesian Estimation and Inference*
 - Gaussian Process

MAP Estimation

- So far, we solve \mathbf{w} by ML estimation:

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

MAP Estimation

- So far, we solve \mathbf{w} by ML estimation:

$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w})$$

- In MAP estimation, we solve

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} P(\mathbf{w} | \mathbb{X}) = \arg \max_{\mathbf{w}} P(\mathbb{X} | \mathbf{w}) P(\mathbf{w})$$

- $P(\mathbf{w})$ models our *preference* or *prior knowledge* about \mathbf{w}

MAP Estimation for Linear Regression

- MAP estimation in linear regression:

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} \log[\mathbf{P}(\mathbb{X} | \mathbf{w}) \mathbf{P}(\mathbf{w})]$$

MAP Estimation for Linear Regression

- MAP estimation in linear regression:

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} \log[\mathbf{P}(\mathbb{X} | \mathbf{w}) \mathbf{P}(\mathbf{w})]$$

- If we assume that $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \beta^{-1} \mathbf{I})$

$$\log[\mathbf{P}(\mathbb{X} | \mathbf{w}) \mathbf{P}(\mathbf{w})] = \log \mathbf{P}(\mathbb{X} | \mathbf{w}) + \log \mathbf{P}(\mathbf{w})$$

MAP Estimation for Linear Regression

- MAP estimation in linear regression:

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} \log[\mathbf{P}(\mathbb{X} | \mathbf{w}) \mathbf{P}(\mathbf{w})]$$

- If we assume that $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \beta^{-1} \mathbf{I})$

$$\begin{aligned} \log[\mathbf{P}(\mathbb{X} | \mathbf{w}) \mathbf{P}(\mathbf{w})] &= \log \mathbf{P}(\mathbb{X} | \mathbf{w}) + \log \mathbf{P}(\mathbf{w}) \propto -\sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 \\ &\quad + \log \sqrt{\frac{1}{(2\pi)^D \det(\beta^{-1} \mathbf{I})}} \exp \left[-\frac{1}{2} (\mathbf{w} - \mathbf{0})^\top (\beta^{-1} \mathbf{I})^{-1} (\mathbf{w} - \mathbf{0}) \right] \end{aligned}$$

MAP Estimation for Linear Regression

- MAP estimation in linear regression:

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} \log[\mathbf{P}(\mathbb{X} | \mathbf{w}) \mathbf{P}(\mathbf{w})]$$

- If we assume that $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \beta^{-1} \mathbf{I})$

$$\begin{aligned} \log[\mathbf{P}(\mathbb{X} | \mathbf{w}) \mathbf{P}(\mathbf{w})] &= \log \mathbf{P}(\mathbb{X} | \mathbf{w}) + \log \mathbf{P}(\mathbf{w}) \propto -\sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 \\ &\quad + \log \sqrt{\frac{1}{(2\pi)^D \det(\beta^{-1} \mathbf{I})}} \exp \left[-\frac{1}{2} (\mathbf{w} - \mathbf{0})^\top (\beta^{-1} \mathbf{I})^{-1} (\mathbf{w} - \mathbf{0}) \right] \\ &\propto -\sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 - \beta \mathbf{w}^\top \mathbf{w} \end{aligned}$$

- $\mathbf{P}(\mathbf{w})$ corresponds to the *weight decay* term in Ridge regression

MAP Estimation for Linear Regression

- MAP estimation in linear regression:

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} \log[\mathbf{P}(\mathbb{X} | \mathbf{w}) \mathbf{P}(\mathbf{w})]$$

- If we assume that $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \beta^{-1} \mathbf{I})$

$$\begin{aligned} \log[\mathbf{P}(\mathbb{X} | \mathbf{w}) \mathbf{P}(\mathbf{w})] &= \log \mathbf{P}(\mathbb{X} | \mathbf{w}) + \log \mathbf{P}(\mathbf{w}) \propto -\sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 \\ &\quad + \log \sqrt{\frac{1}{(2\pi)^D \det(\beta^{-1} \mathbf{I})}} \exp \left[-\frac{1}{2} (\mathbf{w} - \mathbf{0})^\top (\beta^{-1} \mathbf{I})^{-1} (\mathbf{w} - \mathbf{0}) \right] \\ &\propto -\sum_i (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 - \beta \mathbf{w}^\top \mathbf{w} \end{aligned}$$

- $\mathbf{P}(\mathbf{w})$ corresponds to the **weight decay** term in Ridge regression
- MAP estimation provides a way to design complicated yet interpretable regularization terms
 - E.g., we have LASSO by letting $\mathbf{P}(\mathbf{w}) \sim \text{Laplace}(0, b)$ [Proof]
 - We can also let $\mathbf{P}(\mathbf{w})$ be a mixture of Gaussians

Remarks on ML and MAP Estimation

Theorem (Consistency)

*The ML estimator Θ_{ML} is **consistent**, i.e., $\lim_{N \rightarrow \infty} \Theta_{ML} \xrightarrow{\text{Pr}} \Theta^*$ as long as the “true” $P(y|\mathbf{x}; \Theta^*)$ lies within our model \mathbb{F} .*

Remarks on ML and MAP Estimation

Theorem (Consistency)

The ML estimator Θ_{ML} is **consistent**, i.e., $\lim_{N \rightarrow \infty} \Theta_{ML} \xrightarrow{\text{Pr}} \Theta^*$ as long as the “true” $P(y|\mathbf{x}; \Theta^*)$ lies within our model \mathbb{F} .

Theorem (Cramér-Rao Lower Bound [2])

At a fixed (large) number N of examples, no consistent estimator of $\hat{\Theta}$ has a lower expected MSE (mean square error) than the ML estimator Θ_{ML} .

- That is, Θ_{ML} has a low sample complexity (or is statistic efficient)

Remarks on ML and MAP Estimation

Theorem (Consistency)

*The ML estimator Θ_{ML} is **consistent**, i.e., $\lim_{N \rightarrow \infty} \Theta_{ML} \xrightarrow{\text{Pr}} \Theta^*$ as long as the “true” $P(y|\mathbf{x}; \Theta^*)$ lies within our model \mathbb{F} .*

Theorem (Cramér-Rao Lower Bound [2])

At a fixed (large) number N of examples, no consistent estimator of $\hat{\Theta}$ has a lower expected MSE (mean square error) than the ML estimator Θ_{ML} .

- That is, Θ_{ML} has a low sample complexity (or is statistic efficient)
- ML estimation is popular due to its consistency and efficiency

Remarks on ML and MAP Estimation

Theorem (Consistency)

The ML estimator Θ_{ML} is **consistent**, i.e., $\lim_{N \rightarrow \infty} \Theta_{ML} \xrightarrow{\text{Pr}} \Theta^*$ as long as the “true” $P(y|\mathbf{x}; \Theta^*)$ lies within our model \mathbb{F} .

Theorem (Cramér-Rao Lower Bound [2])

At a fixed (large) number N of examples, no consistent estimator of $\hat{\Theta}$ has a lower expected MSE (mean square error) than the ML estimator Θ_{ML} .

- That is, Θ_{ML} has a low sample complexity (or is statistic efficient)
- ML estimation is popular due to its consistency and efficiency
- When N is small that yields overfitting behavior, we can use MAP estimation to **introduce bias** and **reduce variance**

Remarks on ML and MAP Estimation

Theorem (Consistency)

The ML estimator Θ_{ML} is **consistent**, i.e., $\lim_{N \rightarrow \infty} \Theta_{ML} \xrightarrow{\text{Pr}} \Theta^*$ as long as the “true” $P(y|\mathbf{x}; \Theta^*)$ lies within our model \mathbb{F} .

Theorem (Cramér-Rao Lower Bound [2])

At a fixed (large) number N of examples, no consistent estimator of $\hat{\Theta}$ has a lower expected MSE (mean square error) than the ML estimator Θ_{ML} .

- That is, Θ_{ML} has a low sample complexity (or is statistic efficient)
- ML estimation is popular due to its consistency and efficiency
- When N is small that yields overfitting behavior, we can use MAP estimation to **introduce bias** and **reduce variance**

Outline

- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation
- 4 Bayesian Estimation and Inference***
 - Gaussian Process

Bayesian Estimation

- In ML/MAP estimation, we use the estimated $\hat{\Theta}$ as a constant to make prediction:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \hat{\Theta})$$

- $\hat{\Theta} = \arg \max_{\Theta} P(\Theta | \mathbb{X})$

Bayesian Estimation

- In ML/MAP estimation, we use the estimated $\hat{\Theta}$ as a constant to make prediction:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \hat{\Theta})$$

- $\hat{\Theta} = \arg \max_{\Theta} P(\Theta | \mathbb{X})$
- **Bayesian inference** treats Θ as a random variable when making prediction:

$$P(y | \mathbf{x}, \mathbb{X}) = \int_{\Theta} P(y, \Theta | \mathbf{x}, \mathbb{X}) d\Theta = \int P(y | \mathbf{x}, \Theta) P(\Theta | \mathbb{X}) d\Theta$$

Bayesian Estimation

- In ML/MAP estimation, we use the estimated $\hat{\Theta}$ as a constant to make prediction:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}; \hat{\Theta})$$

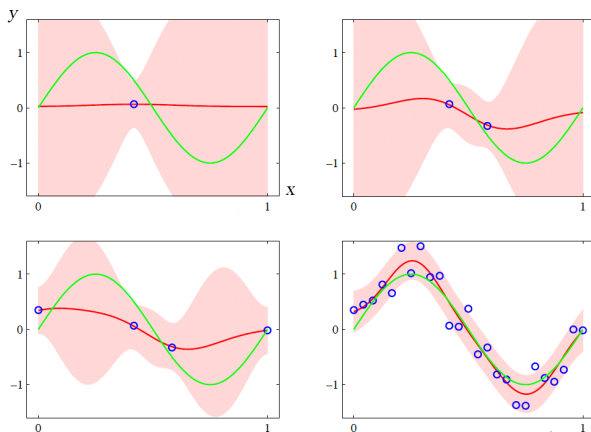
- $\hat{\Theta} = \arg \max_{\Theta} P(\Theta | \mathbb{X})$
- **Bayesian inference** treats Θ as a random variable when making prediction:

$$P(y | \mathbf{x}, \mathbb{X}) = \int_{\Theta} P(y, \Theta | \mathbf{x}, \mathbb{X}) d\Theta = \int P(y | \mathbf{x}, \Theta) P(\Theta | \mathbb{X}) d\Theta$$

- **Entire distribution** $P(y | \mathbf{x}, \mathbb{X})$ is calculated, so we get not only $\hat{y} = \arg \max_y P(y | \mathbf{x}, \mathbb{X})$ but the uncertainty of each prediction
- **Bayesian estimation of Θ** : each prediction considers **all** Θ 's (weighted by their chances $P(\Theta | \mathbb{X})$)

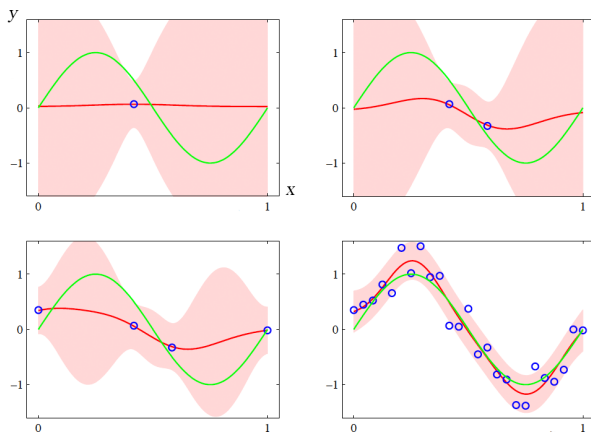
Example: 1D Regression

- Let $y = f^*(x) + \varepsilon$
 - Green line: $f^*(\cdot)$
 - Blue dots: noisy examples



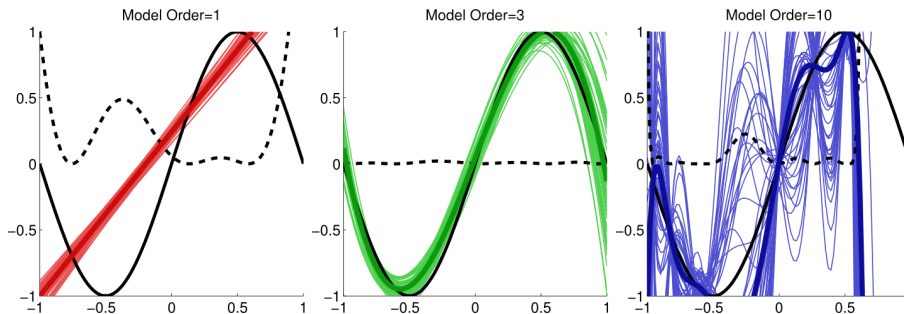
Example: 1D Regression

- Let $y = f^*(x) + \varepsilon$
 - Green line: $f^*(\cdot)$
 - Blue dots: noisy examples
- Red line: predictions by a Bayesian regressor (Gaussian Process)
- Shaded area: confidence intervals of predictions



Bayesian vs. ML Estimation

- Recall the bias-variance trade-off an ML-base polynomial regressor:



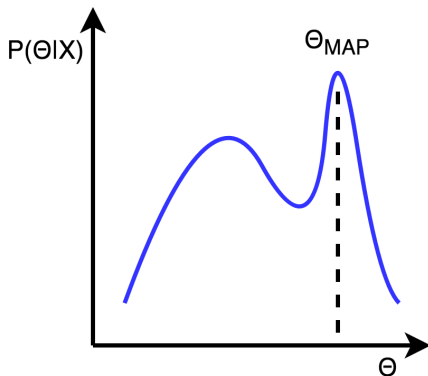
- Bayesian regressor usually generalizes better when the size N of training set is small
 - Avoids high variance $\text{Var}_{\mathbf{X}}(\Theta_{\text{ML}})$

Bayesian vs. MAP Estimation

- MAP gains some benefit of Bayesian approach by incorporating prior as $\text{bias}(\Theta_{\text{MAP}})$
 - Reduces $\text{Var}_{\mathbb{X}}(\Theta_{\text{MAP}})$ when training set is small

Bayesian vs. MAP Estimation

- MAP gains some benefit of Bayesian approach by incorporating prior as bias(Θ_{MAP})
 - Reduces $\text{Var}_{\mathbb{X}}(\Theta_{\text{MAP}})$ when training set is small
- However, does **not** work if Θ_{MAP} is unrepresentative of the majority Θ in $\int P(\mathbf{y}, \Theta | \mathbf{x}, \mathbb{X}) d\Theta = \int P(\mathbf{y} | \mathbf{x}, \Theta) P(\Theta | \mathbb{X}) d\Theta$
- E.g. when $P(\Theta | \mathbb{X})$ is a mixture of Gaussian



Evaluating $P(\mathbf{y} | \mathbf{x}, \mathbb{X})$

$$P(\mathbf{y} | \mathbf{x}, \mathbb{X}) = \int_{\Theta} P(\mathbf{y}, \Theta | \mathbf{x}, \mathbb{X}) d\Theta = \int P(\mathbf{y} | \mathbf{x}, \Theta) P(\Theta | \mathbb{X}) d\Theta$$

- Integral computation make the evaluation challenging
 - The solution may not be tractable in many applications

Evaluating $P(\mathbf{y}|\mathbf{x}, \mathbb{X})$

$$P(\mathbf{y}|\mathbf{x}, \mathbb{X}) = \int_{\Theta} P(\mathbf{y}, \Theta|\mathbf{x}, \mathbb{X})d\Theta = \int P(\mathbf{y}|\mathbf{x}, \Theta)P(\Theta|\mathbb{X})d\Theta$$

- Integral computation make the evaluation challenging
 - The solution may not be tractable in many applications
- Fortunately, in the context of Bayesian linear regression, $P(\mathbf{y}|\mathbf{x}, \mathbb{X})$ can have a simple, closed form [3]

Bayesian Linear Regression

- Assuming that $y = \mathbf{w}^\top \mathbf{x} + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon)$, we have

$$\begin{aligned} P(y|\mathbf{x}, \mathbb{X}) &= \int_{\mathbf{w}} P(y, \mathbf{w} | \mathbf{x}, \mathbb{X}) d\mathbf{w} = \int P(y|\mathbf{x}, \mathbf{w}) P(\mathbf{w} | \mathbb{X}) d\mathbf{w} \\ &= \frac{1}{P(\mathbb{X})} \int P(y|\mathbf{x}, \mathbf{w}) P(\mathbb{X} | \mathbf{w}) P(\mathbf{w}) d\mathbf{w} \\ &= \frac{\prod_{i=1}^N P(\mathbf{x}^{(i)})}{P(\mathbb{X})} \int P(y|\mathbf{x}, \mathbf{w}) \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{w}) d\mathbf{w} \end{aligned}$$

Bayesian Linear Regression

- Assuming that $y = \mathbf{w}^\top \mathbf{x} + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon)$, we have

$$\begin{aligned} P(y|\mathbf{x}, \mathbb{X}) &= \int_{\mathbf{w}} P(y, \mathbf{w} | \mathbf{x}, \mathbb{X}) d\mathbf{w} = \int P(y|\mathbf{x}, \mathbf{w}) P(\mathbf{w} | \mathbb{X}) d\mathbf{w} \\ &= \frac{1}{P(\mathbb{X})} \int P(y|\mathbf{x}, \mathbf{w}) P(\mathbb{X} | \mathbf{w}) P(\mathbf{w}) d\mathbf{w} \\ &= \frac{\prod_{i=1}^N P(\mathbf{x}^{(i)})}{P(\mathbb{X})} \int P(y|\mathbf{x}, \mathbf{w}) \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{w}) d\mathbf{w} \end{aligned}$$

- $P(y|\mathbf{x}, \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}, \sigma_\varepsilon)$ and $P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}^{(i)}, \sigma_\varepsilon), \forall i$

Bayesian Linear Regression

- Assuming that $y = \mathbf{w}^\top \mathbf{x} + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon)$, we have

$$\begin{aligned} P(y|\mathbf{x}, \mathbb{X}) &= \int_{\mathbf{w}} P(y, \mathbf{w} | \mathbf{x}, \mathbb{X}) d\mathbf{w} = \int P(y|\mathbf{x}, \mathbf{w}) P(\mathbf{w} | \mathbb{X}) d\mathbf{w} \\ &= \frac{1}{P(\mathbb{X})} \int P(y|\mathbf{x}, \mathbf{w}) P(\mathbb{X} | \mathbf{w}) P(\mathbf{w}) d\mathbf{w} \\ &= \frac{\prod_{i=1}^N P(\mathbf{x}^{(i)})}{P(\mathbb{X})} \int P(y|\mathbf{x}, \mathbf{w}) \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{w}) d\mathbf{w} \end{aligned}$$

- $P(y|\mathbf{x}, \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}, \sigma_\varepsilon)$ and $P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}^{(i)}, \sigma_\varepsilon)$, $\forall i$
- If we assume that $\mathbf{w} \sim \mathcal{N}$, then $P(y|\mathbf{x}, \mathbf{w}) \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{w})$ can be described by an $N+2$ dimensional Gaussian
 - $(y|\mathbf{x}, \mathbf{w})$, $(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w})$, and \mathbf{w} are (conditionally) independent with each other

Bayesian Linear Regression

- Assuming that $y = \mathbf{w}^\top \mathbf{x} + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon)$, we have

$$\begin{aligned} P(y|\mathbf{x}, \mathbb{X}) &= \int_{\mathbf{w}} P(y, \mathbf{w} | \mathbf{x}, \mathbb{X}) d\mathbf{w} = \int P(y|\mathbf{x}, \mathbf{w}) P(\mathbf{w} | \mathbb{X}) d\mathbf{w} \\ &= \frac{1}{P(\mathbb{X})} \int P(y|\mathbf{x}, \mathbf{w}) P(\mathbb{X} | \mathbf{w}) P(\mathbf{w}) d\mathbf{w} \\ &= \frac{\prod_{i=1}^N P(\mathbf{x}^{(i)})}{P(\mathbb{X})} \int P(y|\mathbf{x}, \mathbf{w}) \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{w}) d\mathbf{w} \end{aligned}$$

- $P(y|\mathbf{x}, \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}, \sigma_\varepsilon)$ and $P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}^{(i)}, \sigma_\varepsilon)$, $\forall i$
- If we assume that $\mathbf{w} \sim \mathcal{N}$, then $P(y|\mathbf{x}, \mathbf{w}) \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{w})$ can be described by an $N+2$ dimensional Gaussian
 - $(y|\mathbf{x}, \mathbf{w})$, $(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w})$, and \mathbf{w} are (conditionally) independent with each other
- Its marginalization, $P(y|\mathbf{x}, \mathbb{X}) = \int P(y, \mathbf{w} | \mathbf{x}, \mathbb{X}) d\mathbf{w}$ is also a Gaussian
 - Gaussian distribution is closed under marginalization

Bayesian Linear Regression

- Assuming that $y = \mathbf{w}^\top \mathbf{x} + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon)$, we have

$$\begin{aligned} P(y|\mathbf{x}, \mathbb{X}) &= \int_{\mathbf{w}} P(y, \mathbf{w} | \mathbf{x}, \mathbb{X}) d\mathbf{w} = \int P(y|\mathbf{x}, \mathbf{w}) P(\mathbf{w} | \mathbb{X}) d\mathbf{w} \\ &= \frac{1}{P(\mathbb{X})} \int P(y|\mathbf{x}, \mathbf{w}) P(\mathbb{X} | \mathbf{w}) P(\mathbf{w}) d\mathbf{w} \\ &= \frac{\prod_{i=1}^N P(\mathbf{x}^{(i)})}{P(\mathbb{X})} \int P(y|\mathbf{x}, \mathbf{w}) \prod_{i=1}^N P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{w}) d\mathbf{w} \end{aligned}$$

- $P(y|\mathbf{x}, \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}, \sigma_\varepsilon)$ and $P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}^{(i)}, \sigma_\varepsilon)$, $\forall i$
- If we assume that $\mathbf{w} \sim \mathcal{N}$, then $P(y|\mathbf{x}, \mathbf{w}) \prod_i P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) P(\mathbf{w})$ can be described by an $N+2$ dimensional Gaussian
 - $(y|\mathbf{x}, \mathbf{w})$, $(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w})$, and \mathbf{w} are (conditionally) independent with each other
- Its marginalization, $P(y|\mathbf{x}, \mathbb{X}) = \int P(y, \mathbf{w} | \mathbf{x}, \mathbb{X}) d\mathbf{w}$ is also a Gaussian
 - Gaussian distribution is closed under marginalization
- Why not model $(y|\mathbf{x}, \mathbb{X}) \sim \mathcal{N}$ in the first place?

Outline

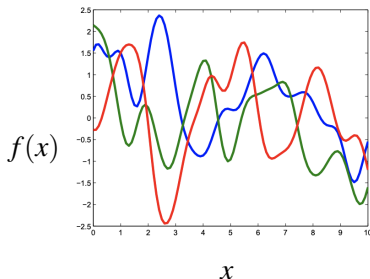
- 1 Probabilistic Models
- 2 Maximum Likelihood Estimation
 - Linear Regression
 - Logistic Regression
- 3 Maximum A Posteriori Estimation
- 4 Bayesian Estimation and Inference***
 - Gaussian Process

Gaussian Process

- Assume a model \mathbb{F} where the domain of each $f(\cdot) \in \mathbb{F}$ consists of only N inputs $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$
- Let $y^{(i)} = f(\mathbf{x}^{(i)}) \in \mathbb{R}$, $\forall i$, we can compactly represent $f(\cdot)$ as a vector $\mathbf{y} = [y^{(1)}, \dots, y^{(N)}]^\top$
- We can specify the probability of $f(\cdot)$ by assuming a distribution over \mathbf{y} , e.g., $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Gaussian Process

- Assume a model \mathbb{F} where the domain of each $f(\cdot) \in \mathbb{F}$ consists of only N inputs $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$
- Let $y^{(i)} = f(\mathbf{x}^{(i)}) \in \mathbb{R}, \forall i$, we can compactly represent $f(\cdot)$ as a vector $\mathbf{y} = [y^{(1)}, \dots, y^{(N)}]^\top$
- We can specify the probability of $f(\cdot)$ by assuming a distribution over \mathbf{y} , e.g., $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- A **stochastic process** is a random distribution over functions in \mathbb{F}
 - Alternatively, it can be a set of random random variables $\{\mathbf{y}^{(i)} \equiv f(\mathbf{x}^{(i)})\}_i$ indexed by time or space i



Gaussian Process (GP)

- A **Gaussian process** is a stochastic process of which the distribution is defined by a mean function $m(\cdot)$ and covariance/**kernel** function $k(\cdot, \cdot)$:

$$\begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} \sim \mathcal{N} \left(\mathbf{m} = \begin{bmatrix} m(\mathbf{x}^{(1)}) \\ \vdots \\ m(\mathbf{x}^{(N)}) \end{bmatrix}, \mathbf{K} = \begin{bmatrix} k(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{bmatrix} \right)$$

- Intuition?

Gaussian Process (GP)

- A **Gaussian process** is a stochastic process of which the distribution is defined by a mean function $m(\cdot)$ and covariance/**kernel** function $k(\cdot, \cdot)$:

$$\begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} \sim \mathcal{N} \left(\mathbf{m} = \begin{bmatrix} m(\mathbf{x}^{(1)}) \\ \vdots \\ m(\mathbf{x}^{(N)}) \end{bmatrix}, \mathbf{K} = \begin{bmatrix} k(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{bmatrix} \right)$$

- Intuition? If $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$ are positively (or negatively) correlated, then $y^{(i)}$ and $y^{(j)}$ should be positively (or negatively) correlated too

Gaussian Process (GP)

- A **Gaussian process** is a stochastic process of which the distribution is defined by a mean function $m(\cdot)$ and covariance/**kernel** function $k(\cdot, \cdot)$:

$$\begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} \sim \mathcal{N} \left(\mathbf{m} = \begin{bmatrix} m(\mathbf{x}^{(1)}) \\ \vdots \\ m(\mathbf{x}^{(N)}) \end{bmatrix}, \mathbf{K} = \begin{bmatrix} k(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{bmatrix} \right)$$

- Intuition? If $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$ are positively (or negatively) correlated, then $y^{(i)}$ and $y^{(j)}$ should be positively (or negatively) correlated too
- Common choices of mean and kernel functions:
 - $m(\cdot) = 0$
 - $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp\left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}{2\tau^2}\right)$ for some fixed $\tau \in \mathbb{R} - \{0\}$
- The kernel matrix \mathbf{K} is usually made positive definite (when $\mathbf{x}^{(i)} \neq \mathbf{x}^{(j)}, \forall i, j$) so it is invertible

Bayesian Regression

- Given N examples $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, how to predict the labels of M unlabeled instances $\mathbb{X}' = \{\mathbf{x}'^{(i)}\}_{i=1}^M$?
- Gaussian process:

$$\begin{bmatrix} \mathbf{y}_N \\ \mathbf{y}_M \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{m}_N \\ \mathbf{m}_M \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{N,N} & \mathbf{K}_{N,M} \\ \mathbf{K}_{M,N} & \mathbf{K}_{M,M} \end{bmatrix}\right),$$

- $\mathbf{m}_N = \mathbf{m}_M = \mathbf{0}$ or $\bar{y}_N \mathbf{1}$, where $\bar{y}_N = \frac{1}{N} \sum_{i=1}^N y^{(i)}$
- \mathbf{y}_M is unknown

Bayesian Regression

- Given N examples $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, how to predict the labels of M unlabeled instances $\mathbb{X}' = \{\mathbf{x}'^{(i)}\}_{i=1}^M$?
- Gaussian process:

$$\begin{bmatrix} \mathbf{y}_N \\ \mathbf{y}_M \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{m}_N \\ \mathbf{m}_M \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{N,N} & \mathbf{K}_{N,M} \\ \mathbf{K}_{M,N} & \mathbf{K}_{M,M} \end{bmatrix}\right),$$

- $\mathbf{m}_N = \mathbf{m}_M = \mathbf{0}$ or $\bar{y}_N \mathbf{1}$, where $\bar{y}_N = \frac{1}{N} \sum_{i=1}^N y^{(i)}$
- \mathbf{y}_M is unknown
- Bayesian inference:

$$P(\mathbf{y}_M | \mathbb{X}', \mathbb{X}) = \mathcal{N}(\mathbf{K}_{M,N} \mathbf{K}_{N,N}^{-1} \mathbf{y}_N, \mathbf{K}_{M,M} - \mathbf{K}_{M,N} \mathbf{K}_{N,N}^{-1} \mathbf{K}_{N,M})$$

- Gaussian distribution is closed under conditioning

Bayesian Regression

- Given N examples $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, how to predict the labels of M unlabeled instances $\mathbb{X}' = \{\mathbf{x}'^{(i)}\}_{i=1}^M$?
- Gaussian process:

$$\begin{bmatrix} \mathbf{y}_N \\ \mathbf{y}_M \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{m}_N \\ \mathbf{m}_M \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{N,N} & \mathbf{K}_{N,M} \\ \mathbf{K}_{M,N} & \mathbf{K}_{M,M} \end{bmatrix}\right),$$

- $\mathbf{m}_N = \mathbf{m}_M = \mathbf{0}$ or $\bar{y}_N \mathbf{1}$, where $\bar{y}_N = \frac{1}{N} \sum_{i=1}^N y^{(i)}$
- \mathbf{y}_M is unknown
- Bayesian inference:

$$P(\mathbf{y}_M | \mathbb{X}', \mathbb{X}) = \mathcal{N}(\mathbf{K}_{M,N} \mathbf{K}_{N,N}^{-1} \mathbf{y}_N, \mathbf{K}_{M,M} - \mathbf{K}_{M,N} \mathbf{K}_{N,N}^{-1} \mathbf{K}_{N,M})$$

- Gaussian distribution is closed under conditioning
- There is **no** explicit training phase
- Predictions: $\hat{\mathbf{y}}_M = \mathbf{K}_{M,N} \mathbf{K}_{N,N}^{-1} \mathbf{y}_N$ (with uncertainty)

Noisy Data

- What if the examples $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ contain noise, i.e., $y = f^*(\mathbf{x}) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$?

Noisy Data

- What if the examples $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ contain noise, i.e., $y = f^*(\mathbf{x}) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$?
- From the i.i.d. noise assumption where $\varepsilon^{(i)}$ and $\varepsilon^{(j)}$ are independent, $\forall i, j$, we have

$$\begin{bmatrix} \mathbf{y}_N \\ \mathbf{y}_M \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{m}_N \\ \mathbf{m}_M \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{N,N} + \sigma_\varepsilon^2 \mathbf{I}_N & \mathbf{K}_{N,M} \\ \mathbf{K}_{M,N} & \mathbf{K}_{M,M} + \sigma_\varepsilon^2 \mathbf{I}_M \end{bmatrix}\right)$$

Noisy Data

- What if the examples $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ contain noise, i.e., $y = f^*(\mathbf{x}) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$?
- From the i.i.d. noise assumption where $\varepsilon^{(i)}$ and $\varepsilon^{(j)}$ are independent, $\forall i, j$, we have

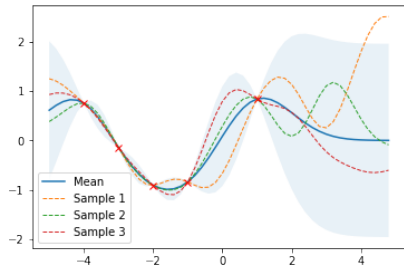
$$\begin{bmatrix} \mathbf{y}_N \\ \mathbf{y}_M \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{m}_N \\ \mathbf{m}_M \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{N,N} + \sigma_\varepsilon^2 \mathbf{I}_N & \mathbf{K}_{N,M} \\ \mathbf{K}_{M,N} & \mathbf{K}_{M,M} + \sigma_\varepsilon^2 \mathbf{I}_M \end{bmatrix}\right)$$

- Bayesian inference:

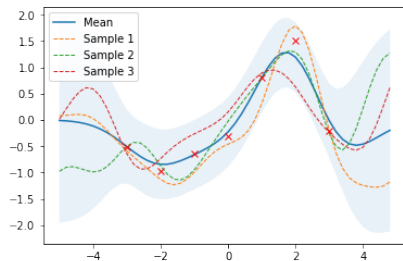
$$\begin{aligned} P(\mathbf{y}_M | \mathbb{X}', \mathbb{X}) &= \mathcal{N}(\mathbf{K}_{M,N}(\mathbf{K}_{N,N} + \sigma_\varepsilon^2 \mathbf{I}_N)^{-1} \mathbf{y}_N, \\ &\quad \mathbf{K}_{M,M} + \sigma_\varepsilon^2 \mathbf{I}_M - \mathbf{K}_{M,N}(\mathbf{K}_{N,N} + \sigma_\varepsilon^2 \mathbf{I}_N)^{-1} \mathbf{K}_{N,M}) \end{aligned}$$

- Predictions: $\hat{\mathbf{y}}_M = \mathbf{K}_{M,N}(\mathbf{K}_{N,N} + \sigma_\varepsilon^2 \mathbf{I}_N)^{-1} \mathbf{y}_N$ (with uncertainty)

Predictions Given Clean and Noisy Data



Clean Data



Noisy Data

Other Choices of Kernels

- Radial basis function (RBF) or exponentiated quadratic kernel:

$$k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp\left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}{2\tau^2}\right)$$

for some fixed $\tau \in \mathbb{R} - \{0\}$

Other Choices of Kernels

- Radial basis function (RBF) or exponentiated quadratic kernel:

$$k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp\left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}{2\tau^2}\right)$$

for some fixed $\tau \in \mathbb{R} - \{0\}$

- Periodic kernel:

$$k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp\left(-\frac{2}{\tau^2} \sin^2\left(\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|}{p} \pi\right)\right)$$

for some fixed $\tau, p \in \mathbb{R} - \{0\}$

- Suitable for periodic data

Other Choices of Kernels

- Radial basis function (RBF) or exponentiated quadratic kernel:

$$k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp\left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}{2\tau^2}\right)$$

for some fixed $\tau \in \mathbb{R} - \{0\}$

- Periodic kernel:

$$k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp\left(-\frac{2}{\tau^2} \sin^2\left(\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|}{p} \pi\right)\right)$$

for some fixed $\tau, p \in \mathbb{R} - \{0\}$

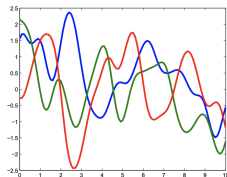
- Suitable for periodic data
- Combined kernel:

$$k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = k^{(1)}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \cdot k^{(2)}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \dots$$

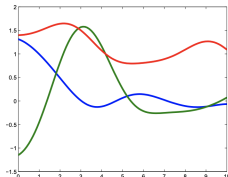
- Has a high value only if all source covariances have a high value (AND operation)

Hyperparameter Tuning

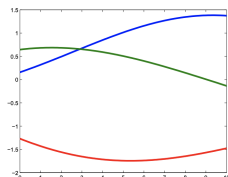
- The τ in the RBF kernel $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}{2\tau^2})$ controls the “smoothness” of the prediction functions



$\tau = 0.5$



$\tau = 2$

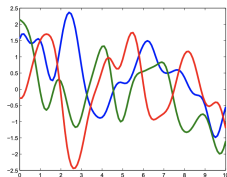


$\tau = 10$

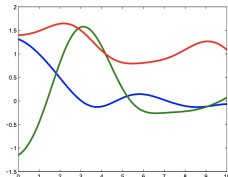
- How to decide the best τ for a given \mathbb{X} ?
 - More generally, how to decide the hyperparameters of chosen kernels?

Hyperparameter Tuning

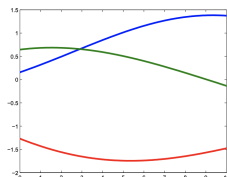
- The τ in the RBF kernel $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2}{2\tau^2})$ controls the “smoothness” of the prediction functions



$\tau = 0.5$



$\tau = 2$



$\tau = 10$

- How to decide the best τ for a given \mathbb{X} ?
 - More generally, how to decide the hyperparameters of chosen kernels?
- We can solve τ using the ML estimation we already familiar with:

$$\begin{aligned}\tau_{\text{ML}} &= \arg \min_{\tau} -\log P(\mathbb{X} | \tau) = \arg \min_{\tau} -\log P(\mathbf{y}_N | \mathbf{X}_N, \tau) \\ &= \arg \min_{\tau} (\mathbf{y}_N - \mathbf{m}_N)^\top \mathbf{K}_{N,N}^{-1} (\mathbf{y}_N - \mathbf{m}_N) + \log \det(\mathbf{K}_{N,N})\end{aligned}$$

- Derivable w.r.t. τ , so can be solved using a gradient-based approach

Parametric vs. Non-Parametric Models

- Probabilistic linear regression and logistic regression are special cases of *parametric models*, whose #parameters is fixed with respect to #data seen
 - $\hat{y} = \mathbf{w}^\top \mathbf{x}$ or $\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x})$
 - *Model complexity grows with data dimension D*

Parametric vs. Non-Parametric Models

- Probabilistic linear regression and logistic regression are special cases of *parametric models*, whose #parameters is fixed with respect to #data seen
 - $\hat{y} = \mathbf{w}^\top \mathbf{x}$ or $\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x})$
 - *Model complexity grows with data dimension D*
- Gaussian process, on the other hand, is a *non-parametric model*
 - $\hat{\mathbf{y}}_M = \mathbf{K}_{M,N} \mathbf{K}_{N,N}^{-1} \mathbf{y}_N$, where each predicted label \hat{y} is a linear combination of the labels in training set
 - *Model complexity grows with N*

Remarks

- Bayesian estimation:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}, \mathbb{X}) = \arg \max_y \int P(y, \Theta | \mathbf{x}, \mathbb{X}) d\Theta$$

- Usually generalizes better given a small training set

Remarks

- Bayesian estimation:

$$\hat{y} = \arg \max_y P(y | \mathbf{x}, \mathbb{X}) = \arg \max_y \int P(y, \Theta | \mathbf{x}, \mathbb{X}) d\Theta$$

- Usually generalizes better given a small training set
- Unfortunately, solution may not be tractable in many applications

Remarks

- Bayesian estimation:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} P(\mathbf{y} | \mathbf{x}, \mathbb{X}) = \arg \max_{\mathbf{y}} \int P(\mathbf{y}, \Theta | \mathbf{x}, \mathbb{X}) d\Theta$$

- Usually generalizes better given a small training set
- Unfortunately, solution may not be tractable in many applications
- Even tractable, incurs high computation cost
 - In GP, each batch of predictions $\hat{\mathbf{y}}_M = \mathbf{K}_{M,N} \mathbf{K}_{N,N}^{-1} \mathbf{y}_N$ may take $O(N^3)$ time
 - Not suitable for large-scale learning tasks

Reference I

- [1] Deepak Roy Chittajallu.
Why is the error function minimized in logistic regression convex?
<http://mathgotchas.blogspot.tw/2011/10/why-is-error-function-minimized-in.html>, 2011.
- [2] Harald Cramér.
Mathematical Methods of Statistics.
Princeton university press, 1946.
- [3] Carl Edward Rasmussen.
Gaussian processes in machine learning.
In *Summer School on Machine Learning*, pages 63–71. Springer, 2003.