

# Lecture 10: Linear Regression

## Introduction to Machine Learning [25737]

Sajjad Amini

Sharif University of Technology

- 1 Approach Definitions
- 2 Linear Regression Model
- 3 Heteroskedastic Regression (Weighted LR)
- 4 Measuring Goodness of Fit
- 5 MAP Estimation (Regularization)
  - Ridge Regression
  - Lasso Regression
- 6 Bayesian Linear Regression

Except explicitly cited, the reference for the material in slides is:

- Murphy, K. P. (2022). *Probabilistic machine learning: an introduction*. MIT press.

# Section 1

## Approach Definitions

## Regression

- Task  $T$ : Finding mapping  $f : \mathbf{x} \mapsto \mathbf{y}$  ( $\mathbf{x} \in \mathcal{X} = \mathbb{R}^D$  and  $y \in \mathbb{R}$ )
- Experience  $E$ : Set of  $N$  input-output pairs  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$
- $P = \frac{1}{N} \sum_{n=1}^N (y_n - f(\mathbf{x}_n; \boldsymbol{\theta}))^2$

## Linear Regression

Similar to classification problems, in regression problems we model  $p(y|\mathbf{x}, \boldsymbol{\theta})$ . Linear regression is the class of regression problem modeling where the expected value of the output is assumed to be a linear function of the input. In other words:

$$\mathbb{E}[y|\mathbf{x}, \boldsymbol{\theta}] = \mathbf{w}^T \mathbf{x}$$

where  $\mathbf{w}$  is a subset of model parameters  $\boldsymbol{\theta}$ .

## Section 2

# Linear Regression Model

# Linear Regression Model

## Linear Regression Model

One model for linear regression can be formulated as:

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|w_0 + \mathbf{w}^T \mathbf{x}, \sigma^2)$$

where:

$\mathbf{w}$	Weights or regression coefficients
$w_0$	Bias or offset
$\sigma^2$	Estimation variance

and  $\boldsymbol{\theta} = [w_0; \mathbf{w}; \sigma^2]$ .

## Vectors Augmentation

Similar to classification models, we usually consider augmented vectors  $[w_0; \mathbf{w}]$  and  $[1; \mathbf{x}]$ , which results in the following model:

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

In this case  $\boldsymbol{\theta} = [\mathbf{w}; \sigma^2]$ .

## Extension to Vector Response $\mathbf{y}$

Consider the situation where response is vector  $\mathbf{y} \in \mathbb{R}^J$  rather than scalar. Then assuming the elements of  $\mathbf{y}$  are independent, we have:

$$p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) = \prod_{j=1}^J \mathcal{N}(y_j | \mathbf{w}_j^T \mathbf{x}, \sigma_j^2)$$

where  $\boldsymbol{\theta} = [\mathbf{w}_1; \dots; \mathbf{w}_J; \sigma_1^2; \dots; \sigma_J^2]$ .

## Feature Transformation

Similar to classification problems, we can use feature transformation to reach a more descriptive models as:

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

where  $\boldsymbol{\theta} = [\mathbf{w}; \sigma^2]$ .



## MLE

Using model formulation, we have:

$$\begin{aligned} p(\mathcal{D}|\boldsymbol{\theta}) &\stackrel{(1)}{=} p(\{y_n\}_{n=1}^N | \{\mathbf{x}_n\}_{n=1}^N, \boldsymbol{\theta}) \stackrel{(2)}{=} \prod_{i=1}^N p(y_n | \mathbf{x}_n, \boldsymbol{\theta}) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_n - \mathbf{w}^T \mathbf{x}_n)^2\right) \end{aligned}$$

where we use mode definition and independence of training samples in equality (1) and (2), respectively. Thus negative log-likelihood is:

$$\begin{aligned} \text{NLL}(\boldsymbol{\theta}) &= -\sum_{n=1}^N \log \left[ \left( \frac{1}{2\pi\sigma^2} \right) \exp \left( -\frac{1}{2\sigma^2} (y_n - \mathbf{w}^T \mathbf{x}_n)^2 \right) \right] \\ &= \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \hat{y}_n)^2 + \frac{N}{2} \log(2\pi\sigma^2) \end{aligned}$$

where we define  $\hat{y}_n \triangleq \mathbf{w}^T \mathbf{x}_n$ .

## Converting Summation into Matrix Form

We can easily show that:

$$\text{RSS}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2 = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

where  $\text{RSS}(\cdot)$  stand for *residual sum of squares* function and:

$$\mathbf{X} = \begin{bmatrix} - & \mathbf{x}_1 & - \\ & \vdots & \\ - & \mathbf{x}_n & - \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

Thus the negative log-likelihood can be written as:

$$\text{NLL}(\boldsymbol{\theta}) = \frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \frac{N}{2} \log(2\pi\sigma^2)$$

## MLE

The optimal value for  $\mathbf{w}$  and  $\sigma^2$  can be calculated as:

$$\begin{aligned}\nabla_{\mathbf{w}} \text{NLL}(\boldsymbol{\theta}) = \mathbf{0} &\Rightarrow \nabla_{\mathbf{w}} \text{RSS}(\mathbf{w}) = \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = \mathbf{0} \\ \Rightarrow \hat{\mathbf{w}}_{mle} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ \nabla_{\sigma} \text{NLL}(\boldsymbol{\theta}) = 0 &\Rightarrow \hat{\sigma}_{mle}^2 = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{\mathbf{w}}_{mle}^T \mathbf{x}_n)^2 = \frac{2}{N} \text{RSS}(\hat{\mathbf{w}}_{mle})\end{aligned}$$

Note that  $\hat{\mathbf{w}}_{mle} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  is known as *Ordinary Least Squares* (OLS) solution.

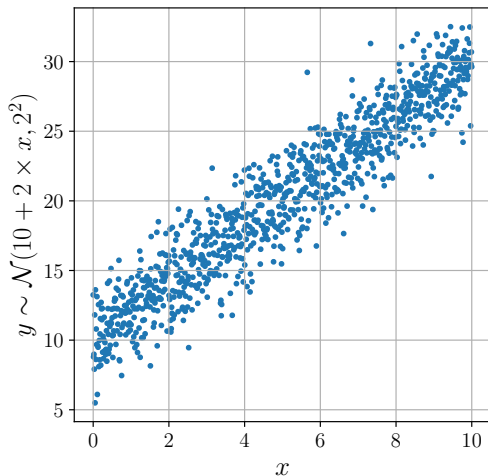
# Sample: Modeling a Device Measurement

## Data Generation

Assume we have device that generates its measurements based on the following rule:

$$\begin{cases} y_i \sim \mathcal{N}(\mu_i, \sigma^2) \\ \mu_i = 10 + 2 \times x_i \\ \sigma^2 = 2^2 \end{cases}$$

On the right, we see a realization of this data for  $N = 1000$ .



# Sample: Modeling a Device Measurement

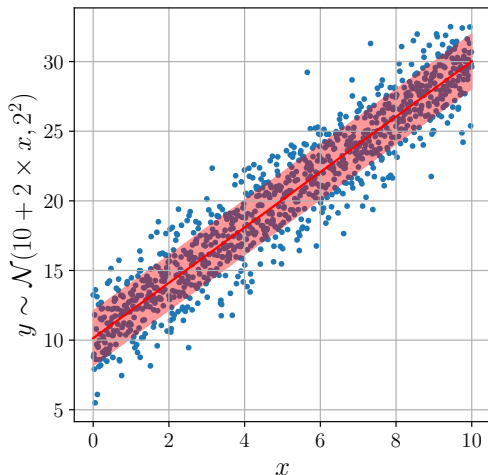
## Data Generation

*Solution:* Based on the MLE formulation, the solution is:

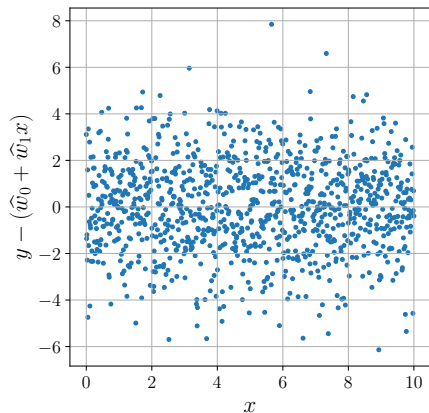
$$\begin{aligned}\hat{\mathbf{w}}_{mle} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= [10.14, 1.99]^T\end{aligned}$$

$$\begin{aligned}\hat{\sigma}_{mle}^2 &= \frac{1}{1000} \sum_{n=1}^{1000} (y_n - \hat{\mathbf{w}}_{mle}^T \mathbf{x}_n)^2 \\ &= 3.85\end{aligned}$$

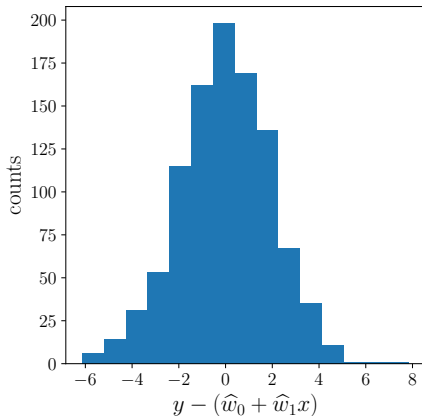
On the right, we see the solution.



# Sample: Modeling a Device Measurement



(a) Difference



(b) Difference Histogram

Figure: Visualization of difference

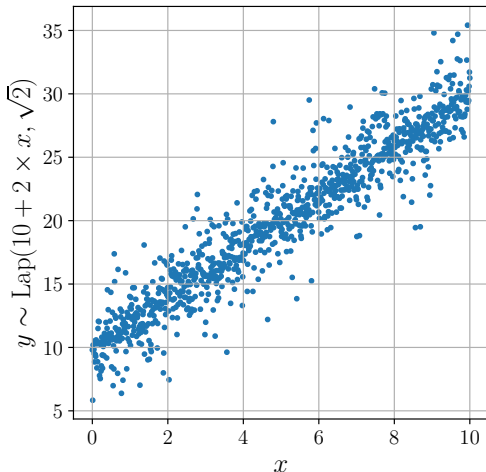
# Sample: Modeling a Device Measurement

## Data Generation

Assume we have device that generates its measurements based on the following rule:

$$\begin{cases} y_i \sim \text{Lap}(\mu_i, b) \\ \mu_i = 10 + 2 \times x_i \\ b = \sqrt{2} \end{cases}$$

On the right, we see a realization of this data for  $N = 1000$ .



# Sample: Modeling a Device Measurement

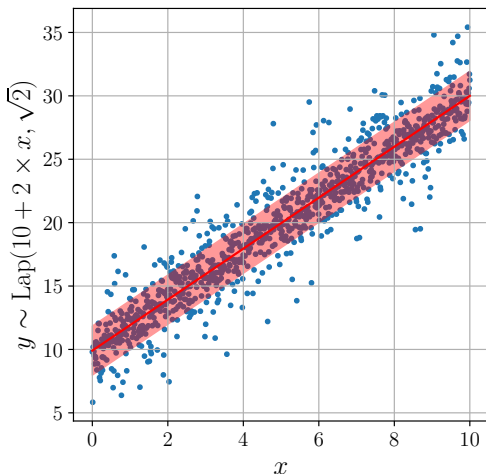
## Data Generation

*Solution:* Based on the MLE formulation, the solution is:

$$\begin{aligned}\hat{\mathbf{w}}_{mle} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= [9.91, 2.01]^T\end{aligned}$$

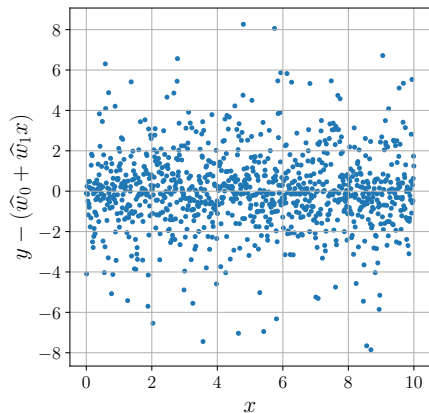
$$\begin{aligned}\hat{\sigma}_{mle}^2 &= \frac{1}{1000} \sum_{n=1}^{1000} (y_n - \hat{\mathbf{w}}_{mle}^T \mathbf{x}_n)^2 \\ &= 3.74\end{aligned}$$

On the right, we see the solution.

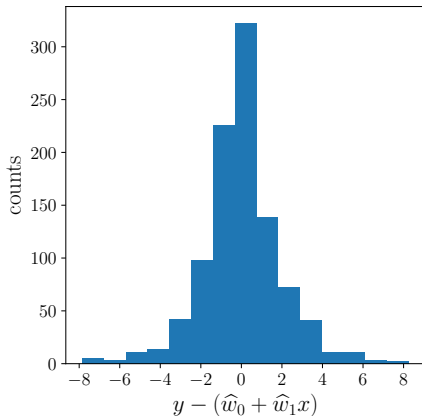




# Sample: Modeling a Device Measurement



(a) Difference



(b) Difference Histogram

Figure: Visualization of difference

## Section 3

# Heteroskedastic Regression (Weighted LR)

## Heteroskedastic Regression

Heteroskedastic Regression assume the following model for data:

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|\mathbf{w}^T \mathbf{x}, \sigma^2(\mathbf{x})) = \frac{1}{\sqrt{2\pi\sigma^2(\mathbf{x})}} \exp\left(-\frac{1}{2\sigma^2(\mathbf{x})}(y - \mathbf{w}^T \mathbf{x})^2\right)$$

## MLE

For the simplicity, assume we have access to  $\{\sigma^2(\mathbf{x}_n)\}_{n=1}^N$ , then:

$$p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2(\mathbf{x}_n)}} \exp\left(-\frac{1}{2\sigma^2(\mathbf{x}_n)}(y_n - \mathbf{w}^T \mathbf{x}_n)^2\right)$$
$$\mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \boldsymbol{\Lambda}^{-1}), \quad \boldsymbol{\Lambda} = \text{diag}\left(\frac{1}{\sigma^2(\mathbf{x}_n)}\right)$$

## MLE

$$p(\mathcal{D}|\boldsymbol{\theta}) = \frac{1}{|2\pi\boldsymbol{\Lambda}^{-1}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T \boldsymbol{\Lambda}(\mathbf{y} - \mathbf{X}\mathbf{w})\right)$$

$$\Rightarrow \text{NLL}(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T \boldsymbol{\Lambda}(\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$\Rightarrow \nabla_{\mathbf{w}} \text{NLL}(\boldsymbol{\theta}) = \mathbf{0} \Rightarrow \mathbf{X}^T \boldsymbol{\Lambda} \mathbf{X} \mathbf{w} - \mathbf{X}^T \boldsymbol{\Lambda} \mathbf{y} = \mathbf{0} \Rightarrow \hat{\mathbf{w}}_{mle} = (\mathbf{X}^T \boldsymbol{\Lambda} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Lambda} \mathbf{y}$$

The above is known as *Weighted Least Squares* (WLS) estimate.

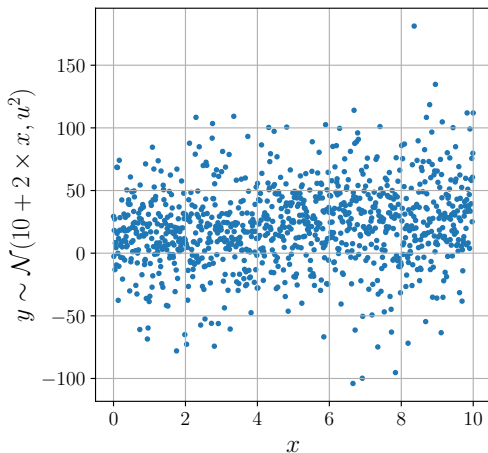
# Sample: Modeling a Device Measurement

## Data Generation

Assume we have device that generates its measurements based on the following rule:

$$\begin{cases} y_i \sim \mathcal{N}(\mu_i, \sigma_i^2) \\ \mu_i = \beta_0 + \beta_1 x_i \\ \sigma_i^2 = u_i^2 \\ u_i \propto U(10, 50) \end{cases}$$

On the right, we see a realization of this data for  $N = 1000$ .



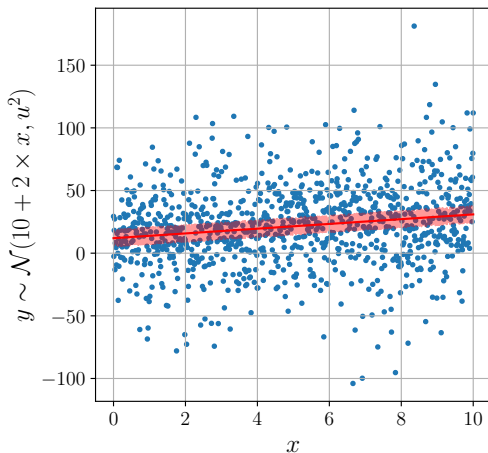
# Sample: Modeling a Device Measurement

## MLE Solution

*Solution:* Based on OLS, we have:

$$\begin{aligned}\hat{\mathbf{w}}_{mle} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= [11.97, 1.90]^T\end{aligned}$$

Below we see the solution.



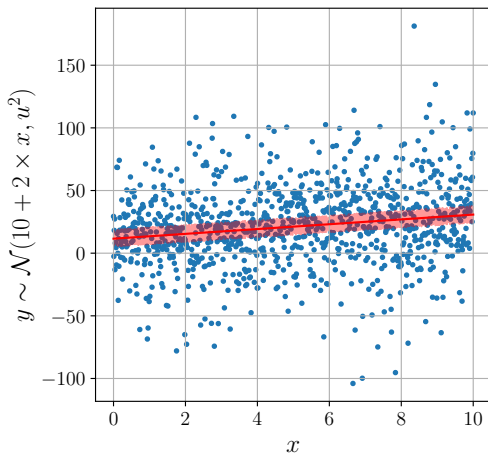
# Sample: Modeling a Device Measurement

## MLE Solution

*Solution:* Based on WLS, we have:

$$\begin{aligned}\hat{\mathbf{w}}_{mle} &= (\mathbf{X}^T \mathbf{\Lambda} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Lambda} \mathbf{y} \\ &= [11.57, 1.92]^T\end{aligned}$$

Below we see the solution.



## Section 4

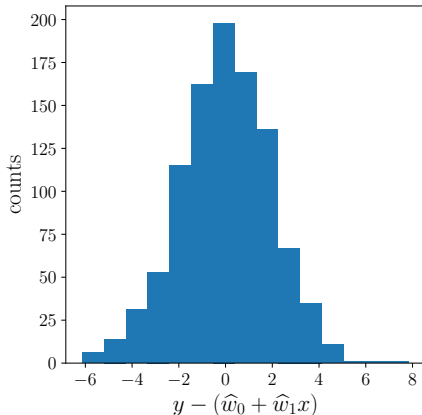
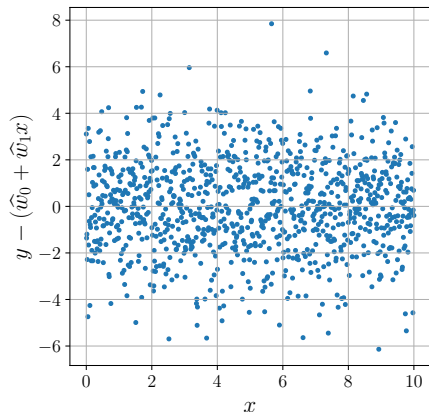
# Measuring Goodness of Fit



# Residual Plots

## Residual Plots

For one dimensional inputs, we can plot the residual  $r_n = y_n - \hat{y}_n$  vs the input  $x_n$ . The resulting plot is called residual plot. This plot should be similar to samples of  $\mathcal{N}(0, \sigma^2)$ .



## Coefficient of Determination

Coefficient of Determination is defined as:

$$R^2 \triangleq 1 - \frac{\sum_{n=1}^N (\hat{y}_n - y_n)^2}{\sum_{n=1}^N (\bar{y} - y_n)^2} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

where  $\bar{y} = \frac{1}{N} \sum_{n=1}^N y_n$  and  $\text{TSS} = \sum_{n=1}^N (\bar{y}_n - y_n)^2$  stands for total sum of squares (TSS). You can show that  $0 \leq R^2 \leq 1$

## Section 5

# MAP Estimation (Regularization)

## Subsection 1

# Ridge Regression

## Ridge Regression

To avoid overfitting in linear regression, similar to classification models, we can assume the weight vector to come from the following prior:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \tau^2 \mathbf{I})$$

where  $\lambda$  is a hyper-parameter. The resulting MAP estimation is known as Ridge Regression and is formulated as:

$$\hat{\mathbf{w}}_{map} = \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathbf{w}|\mathcal{D})p(\mathbf{w})$$

## MAP

Assuming the  $\sigma$  to be known, we have:

$$\begin{aligned} p(\mathbf{w}|\mathcal{D}) &\propto \overbrace{\prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_n - \mathbf{w}^T \mathbf{x}_n)^2\right)}^{p(\mathcal{D}|\mathbf{w})} \overbrace{\frac{1}{|2\pi\tau^2\mathbf{I}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}\mathbf{w}^T(\tau^2\mathbf{I})^{-1}\mathbf{w}\right)}^{p(\mathbf{w})} \\ &\propto \frac{1}{|2\pi\sigma^2\mathbf{I}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\sigma^2\mathbf{I})^{-1}(\mathbf{y} - \mathbf{X}\mathbf{w})^T\right) \\ &\quad \frac{1}{|2\pi\tau^2\mathbf{I}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}\mathbf{w}^T(\tau^2\mathbf{I})^{-1}\mathbf{w}\right) \end{aligned}$$

Thus we have:

$$\begin{aligned} \operatorname{argmax}_{\mathbf{w}} p(\mathbf{w}|\mathcal{D}) &\equiv \operatorname{argmin}_{\mathbf{w}} -\log p(\mathbf{w}|\mathcal{D}) = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \frac{1}{2\tau^2} \|\mathbf{w}\|^2 \\ &\equiv \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \overbrace{\frac{\sigma^2}{\tau^2}}^{\lambda} \|\mathbf{w}\|^2 \end{aligned}$$

## MAP

$$J(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda\|\mathbf{w}\|^2$$

$$\Rightarrow \nabla_{\mathbf{w}} J(\mathbf{w}) = 2(\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} + \lambda \mathbf{w}) = \mathbf{0}$$

$$\Rightarrow \hat{\mathbf{w}}_{map} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} = \left( \sum_n \mathbf{x}_n \mathbf{x}_n^T + \lambda \mathbf{I} \right)^{-1} \left( \sum_n y_n \mathbf{x}_n \right)$$

## Subsection 2

# Lasso Regression



## Feature Selection

Assume the problem we encountered in ridge regression:

$$\operatorname{argmin}_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

Now change the  $\ell_2$  norm with  $\ell_0$  norm as:

$$\operatorname{argmin}_{\mathbf{w}} \underbrace{\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_0}_{J_0(\mathbf{w})}$$

where  $\|\mathbf{w}\|_0 = \sum_{d=1}^D \mathbb{I}(|w_d| > 0)$ . In the above formulation, the problem is solved by using reduced number of features.

# Sample: Modeling a Device Measurement

## Data Generation

Assume we have device that generates its measurements based on  $\begin{cases} y_i \sim \mathcal{N}(\mu_i, \sigma^2) \\ \mu_i = 10 + 2 \times x_i \\ \sigma^2 = 2^2 \end{cases}$

Then we have the following states:

- $\beta_0 \neq 0, \beta_1 \neq 0 \Rightarrow J(\tilde{\mathbf{w}}) = 3848.2 + \lambda \times 2$
- $\beta_0 \neq 0, \beta_1 = 0 \Rightarrow J(\tilde{\mathbf{w}}) = 36842.5 + \lambda \times 1$
- $\beta_0 = 0, \beta_1 \neq 0 \Rightarrow J(\tilde{\mathbf{w}}) = 29584.4 + \lambda \times 1$
- $\beta_0 = 0, \beta_1 = 0 \Rightarrow J(\tilde{\mathbf{w}}) = 439953.5 + \lambda \times 0$

Thus we have the following result for the  $J_0(\mathbf{w})$  problem:

$$\hat{\mathbf{w}} = \begin{cases} [0, 0]^T & (J(\hat{\mathbf{w}}) = 439953.5) & 439953.5 - 29584.4 < \lambda \\ [0, 3.51]^T & (J(\hat{\mathbf{w}}) = 29584.4 + \lambda \times 1) & 29584.4 - 3848.2 < \lambda < 439953.5 - 29584.4 \\ [10.13, 1.98]^T & (J(\hat{\mathbf{w}}) = 3848.2 + \lambda \times 2) & \lambda < 29584.4 - 3848.2 \end{cases}$$

# Sample: Modeling a Device Measurement

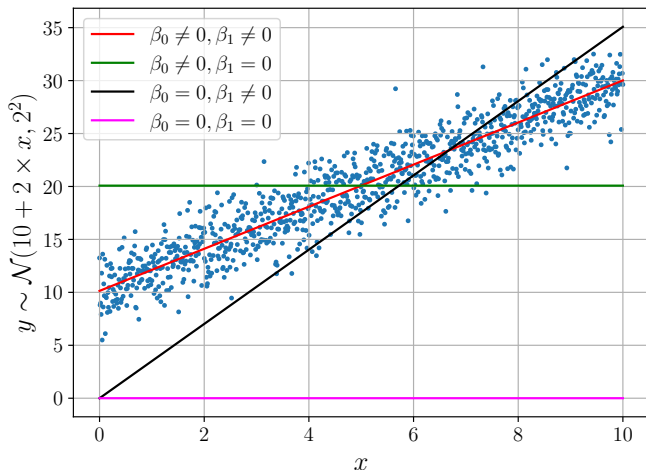


Figure: Realization of device data for  $N = 1000$

## Lasso Regression

Assume dimensions of weight vector are independently and identically distributed as  $\text{Lap}(w|0, b)$ . Then the prior distribution over weight vector is:

$$p(\mathbf{w}) = \prod_{d=1}^D \frac{1}{2b} \exp\left(-\frac{|w_d|}{b}\right)$$

where  $b$  is a hyper-parameter. The resulting MAP estimation is known as Lasso (Least Absolute Shrinkage and Selection Operator) Regression and is formulated as:

$$\hat{\mathbf{w}}_{map} = \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathbf{w}|\mathcal{D})p(\mathbf{w})$$

## MAP

Assuming the  $\sigma$  to be known, we have:

$$\begin{aligned} p(\mathbf{w}|\mathcal{D}) &\propto \overbrace{\prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_n - \mathbf{w}^T \mathbf{x}_n)^2\right)}^{p(\mathcal{D}|\mathbf{w})} \overbrace{\frac{1}{(2b)^D} \exp\left(-\frac{1}{b} \sum_{d=1}^D |w_d|\right)}^{p(\mathbf{w})} \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T (\sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \mathbf{X}\mathbf{w})\right) \exp\left(-\frac{1}{b} \|\mathbf{w}\|_1\right) \end{aligned}$$

Thus we have:

$$\begin{aligned} \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathbf{w}|\mathcal{D}) &\equiv \underset{\mathbf{w}}{\operatorname{argmin}} -\log p(\mathbf{w}|\mathcal{D}) = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \frac{1}{b} \|\mathbf{w}\|_1 \\ &\equiv \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \overbrace{\frac{\lambda}{2\sigma^2}} \|\mathbf{w}\|_1 \end{aligned}$$

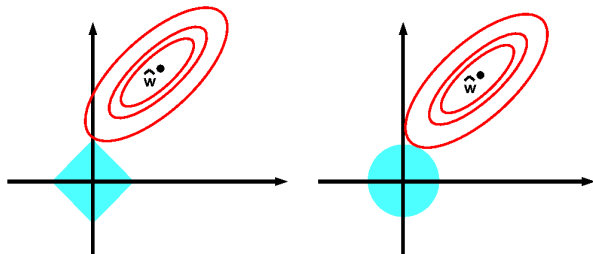
# Connection to Feature Selection

## Lagrangian Interpretation

Using Lagrangian interpretation, we have:

$$\text{RidgeRegression} : \min_{\mathbf{w}} \text{NLL}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \equiv \begin{cases} \min_{\mathbf{w}} \text{NLL}(\mathbf{w}) \\ \text{s.t. } \|\mathbf{w}\|_2^2 \leq B \end{cases}$$

$$\text{LassoRegression} : \min_{\mathbf{w}} \text{NLL}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 \equiv \begin{cases} \min_{\mathbf{w}} \text{NLL}(\mathbf{w}) \\ \text{s.t. } \|\mathbf{w}\|_1 \leq C \end{cases}$$



## MAP

For the solution of Lasso regression, we only consider the case where  $\mathbf{X}^T \mathbf{X} = \mathbf{I}$ . In this case, we have:

$$\hat{\mathbf{w}}_{mle} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{y}$$

If we use the above, we have:

$$\begin{aligned} J(\mathbf{w}) &= \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|_1 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \|\mathbf{w}\|_1 \\ &= \mathbf{y}^T \mathbf{y} - 2 \mathbf{y}^T \underbrace{\mathbf{X} \mathbf{w}}_{\hat{\mathbf{w}}_{mle}^T} + \mathbf{w}^T \underbrace{\mathbf{X}^T \mathbf{X}}_{\mathbf{I}} \mathbf{w} + \lambda \|\mathbf{w}\|_1 + \hat{\mathbf{w}}_{mle}^T \hat{\mathbf{w}}_{mle} - \underbrace{\hat{\mathbf{w}}_{mle}^T \hat{\mathbf{w}}_{mle}}_{\mathbf{y}^T \mathbf{X} \mathbf{X}^T \mathbf{y}} \\ &= \|\mathbf{w} - \hat{\mathbf{w}}_{mle}\|^2 + \lambda \|\mathbf{w}\|_1 + \mathbf{y}^T (\mathbf{I} - \mathbf{X} \mathbf{X}^T) \mathbf{y} \end{aligned}$$

Thus we have:

$$\begin{aligned} \min_{\mathbf{w}} J(\mathbf{w}) &\equiv \min_{\mathbf{w}} \|\mathbf{w} - \hat{\mathbf{w}}_{mle}\|^2 + \lambda \|\mathbf{w}\|_1 \\ &\equiv \min_{w_i} (w_i - \hat{w}_{(mle)_i})^2 + \lambda |w_i|, \quad i = 1, \dots, D \end{aligned}$$

## MAP

We should solve the problem of the following form:

$$\hat{w}_{lasso} = \underset{w}{\operatorname{argmin}} \overbrace{(w - \hat{w}_{mle})^2 + \lambda|w|}^{J(w)}$$

Now assume two cases:

- $\hat{w}_{mle} > 0 \Rightarrow w > 0$ , then:

$$J(w) = (w - \hat{w}_{mle})^2 + \lambda w \Rightarrow \frac{d}{dw} J(w) = 2(w - \hat{w}_{mle}) + \lambda = 0 \Rightarrow w = \hat{w}_{mle} - \frac{\lambda}{2}$$

- $\hat{w}_{mle} \leq 0 \Rightarrow w \leq 0$ , then:

$$J(w) = (w - \hat{w}_{mle})^2 - \lambda w \Rightarrow \frac{d}{dw} J(w) = 2(w - \hat{w}_{mle}) - \lambda = 0 \Rightarrow w = \hat{w}_{mle} + \frac{\lambda}{2}$$

$$\text{Altogether we have: } \hat{w}_{lasso} = \begin{cases} \max\{\hat{w}_{mle} - \frac{\lambda}{2}, 0\} & \hat{w}_{mle} > 0 \\ \min\{\hat{w}_{mle} + \frac{\lambda}{2}, 0\} & \hat{w}_{mle} \leq 0 \end{cases} = \mathcal{S}(\hat{w}_{mle}, \frac{\lambda}{2})$$



# Soft Thresholding Operator

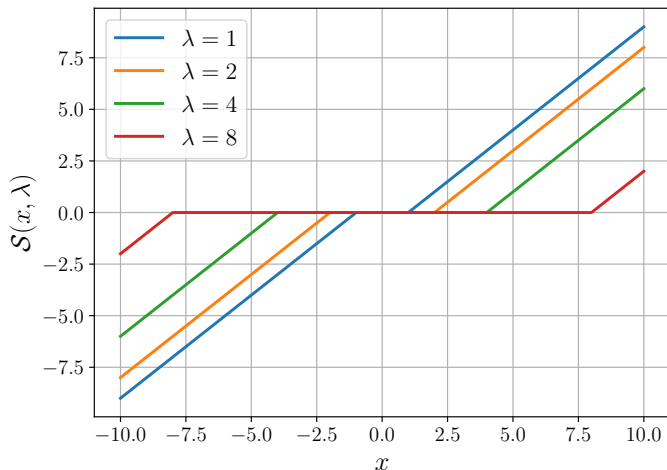


Figure: Soft Thresholding Operator Curve

## Section 6

# Bayesian Linear Regression

# Bayesian Linear Regression

## Prior and Likelihood

Assume the following prior distribution:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \check{\mathbf{w}}, \check{\Sigma})$$

On the other hand, we see before that the likelihood can be written as (we assume  $\sigma^2$  to be known):

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

## Bayes Rule for Gaussian

If  $\begin{cases} p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z) \\ p(\mathbf{y} | \mathbf{z}) = \mathcal{N}(\mathbf{y} | \mathbf{W}\mathbf{z} + \mathbf{b}, \boldsymbol{\Sigma}_y) \end{cases}$ , then:

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

## Posterior

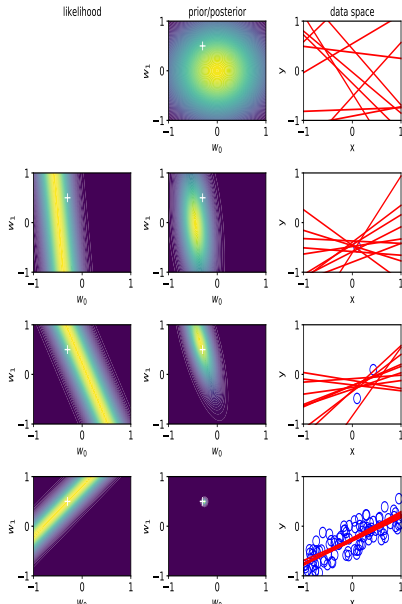
Thus the posterior can be calculated as:

$$p(\mathbf{w}|\mathcal{D}) \propto \mathcal{N}(\mathbf{w}|\check{\mathbf{w}}, \check{\Sigma})\mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \sigma^2\mathbf{I}) = \mathcal{N}(\mathbf{w}|\hat{\mathbf{w}}, \hat{\Sigma})$$

where we have:

$$\hat{\mathbf{w}} \triangleq \hat{\Sigma}(\check{\Sigma}^{-1}\check{\mathbf{w}} + \frac{1}{\sigma^2}\mathbf{X}^T\mathbf{y})$$
$$\hat{\Sigma} \triangleq (\check{\Sigma}^{-1} + \frac{1}{\sigma^2}\mathbf{X}^T\mathbf{X})^{-1}$$

# Bayesian Linear Regression



## Normalization Constant in Bayes' Rule for Gaussian

In Bayes' rule for Gaussian, we have:

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

where the normalization factor is:

$$p(\mathbf{y}) = \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)\mathcal{N}(\mathbf{y}|\mathbf{W}\mathbf{z} + \mathbf{b}, \boldsymbol{\Sigma}_y)d\mathbf{z} = \mathcal{N}(\mathbf{y}|\mathbf{W}\boldsymbol{\mu}_z + \mathbf{b}, \boldsymbol{\Sigma}_y + \mathbf{W}\boldsymbol{\Sigma}_z\mathbf{W}^T)$$

## Computing Posterior Prediction

$$p(y|\mathbf{x}, \mathcal{D}) = \int \mathcal{N}(\mathbf{w}|\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})\mathcal{N}(y|\mathbf{x}^T\mathbf{w}, \sigma^2)d\mathbf{w} = \mathcal{N}(y|\hat{\boldsymbol{\mu}}^T\mathbf{x}, \hat{\sigma}^2(\mathbf{x}))$$
$$\hat{\sigma}^2(\mathbf{x}) = \sigma^2 + \mathbf{x}^T\hat{\boldsymbol{\Sigma}}\mathbf{x}$$