

Lec 11 Machine Learning

11-2 Linear Model

Yang Shu

School of Data Science and Engineering

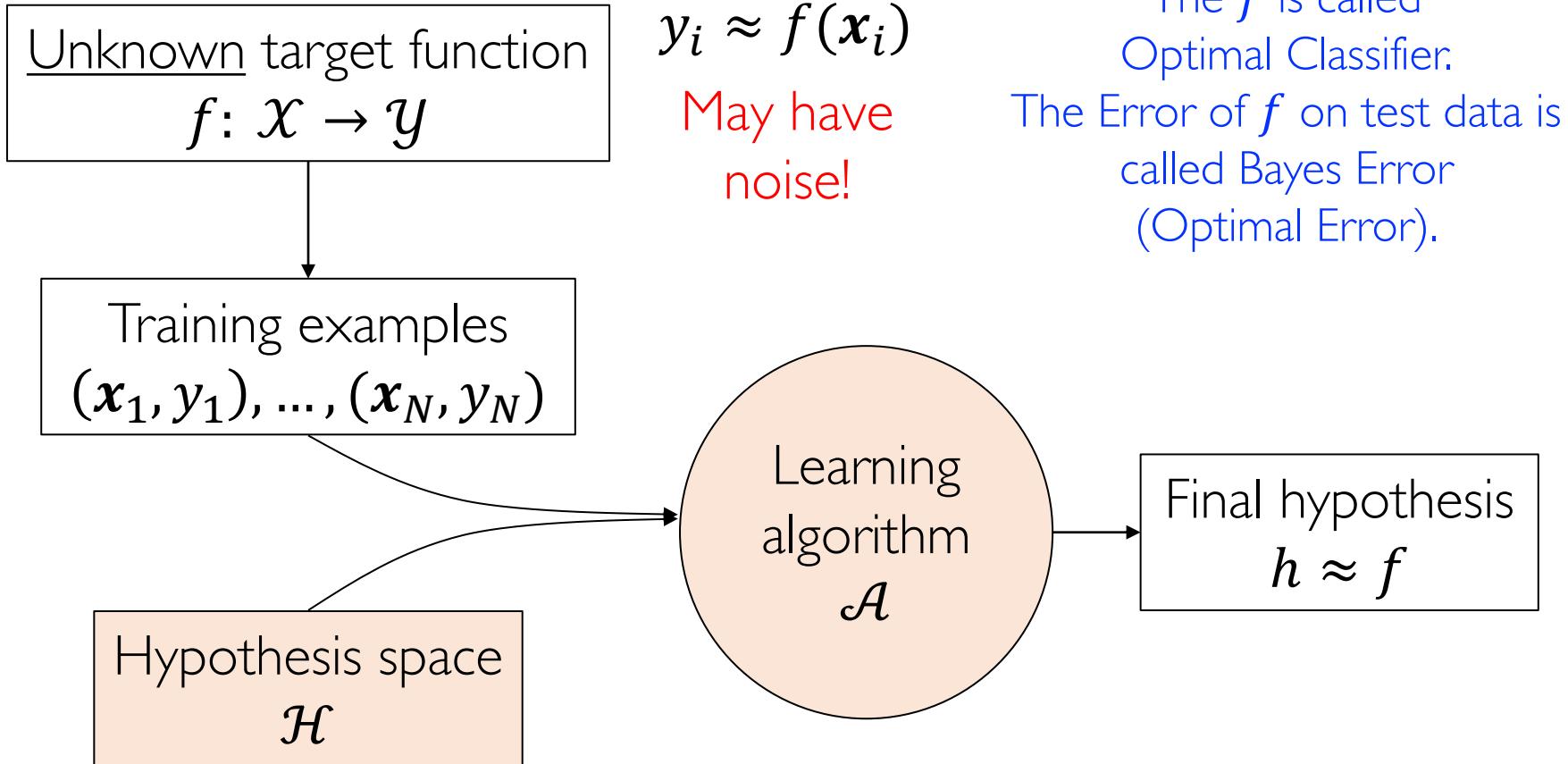
yshu@dase.ecnu.edu.cn

[Acknowledgement: Slides are adapted from Machine Learning Course, Mingsheng Long, THU]

Lec 11-2 Linear Model

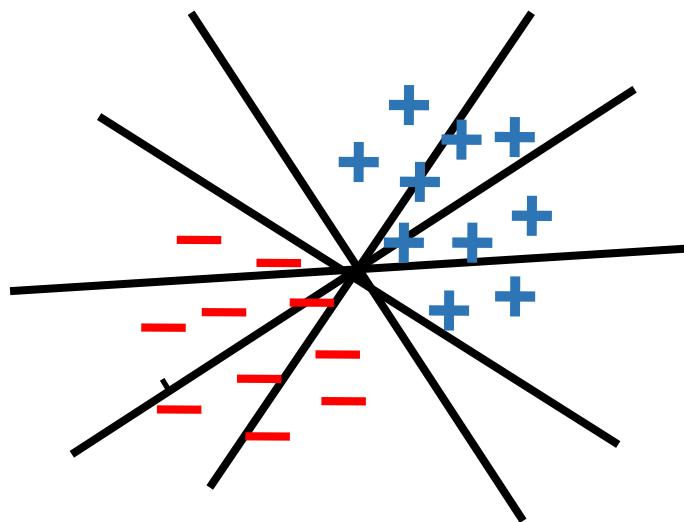
- Linear Regression
 - Optimization
 - Statistical View: Maximum Likelihood Estimation
- Nonlinearization
- Regularization
 - L1 and L2 Norm
 - Statistical View: Maximum a Posteriori Estimation
- Linear Classification
 - Logistic Regression
 - Softmax Regression

Components of Learning



Hypothesis Space

- A **hypothesis space** \mathcal{H} is a set of functions that maps $\mathcal{X} \rightarrow \mathcal{Y}$.
 - It is the collection of prediction functions we are **choosing** from.
- We want hypothesis space that...
 - Includes only those functions that have desired **regularity** (正则性)
 - Continuity (连续性)
 - Smoothness (光滑性)
 - Simplicity (简单性)
 - **Easy** to work with
- An example hypothesis space:
 - All linear hyperplanes for classification



How to find the best?

Loss Function

- **Loss function:** $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ measures the difference between a prediction $h(\mathbf{x})$ and an actual output y :

$$\ell(y, h(\mathbf{x})) = (y - h(\mathbf{x}))^2 \text{ (Squared loss for Regression)}$$

$$\ell(y, h(\mathbf{x})) = \mathbf{1}[y \neq h(\mathbf{x})] \text{ (01-loss for Classification)}$$

- The canonical training procedure of machine learning:

Fit dataset with
best hypothesis

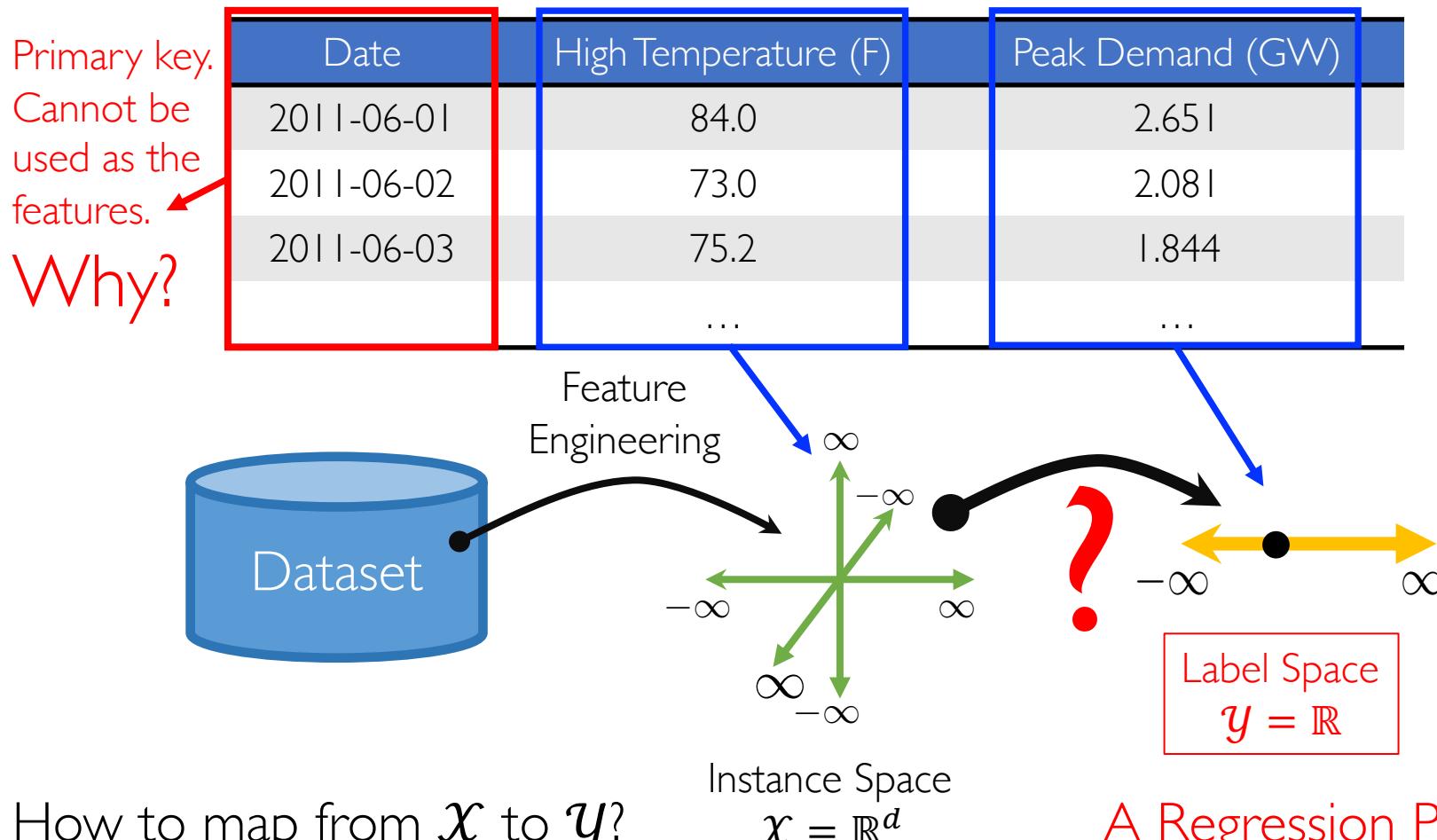
$$\hat{\epsilon}(h) = \min_{\theta} \sum_{i=1}^m \ell(h_{\theta}(\mathbf{x}_i), y_i)$$

θ is parameters of
the hypothesis h

- Virtually every machine learning algorithm has this form, just specify
 - What is the **hypothesis function**?
 - What is the **loss function**?
 - How do we solve the **training problem**?

An Example: Electricity Prediction

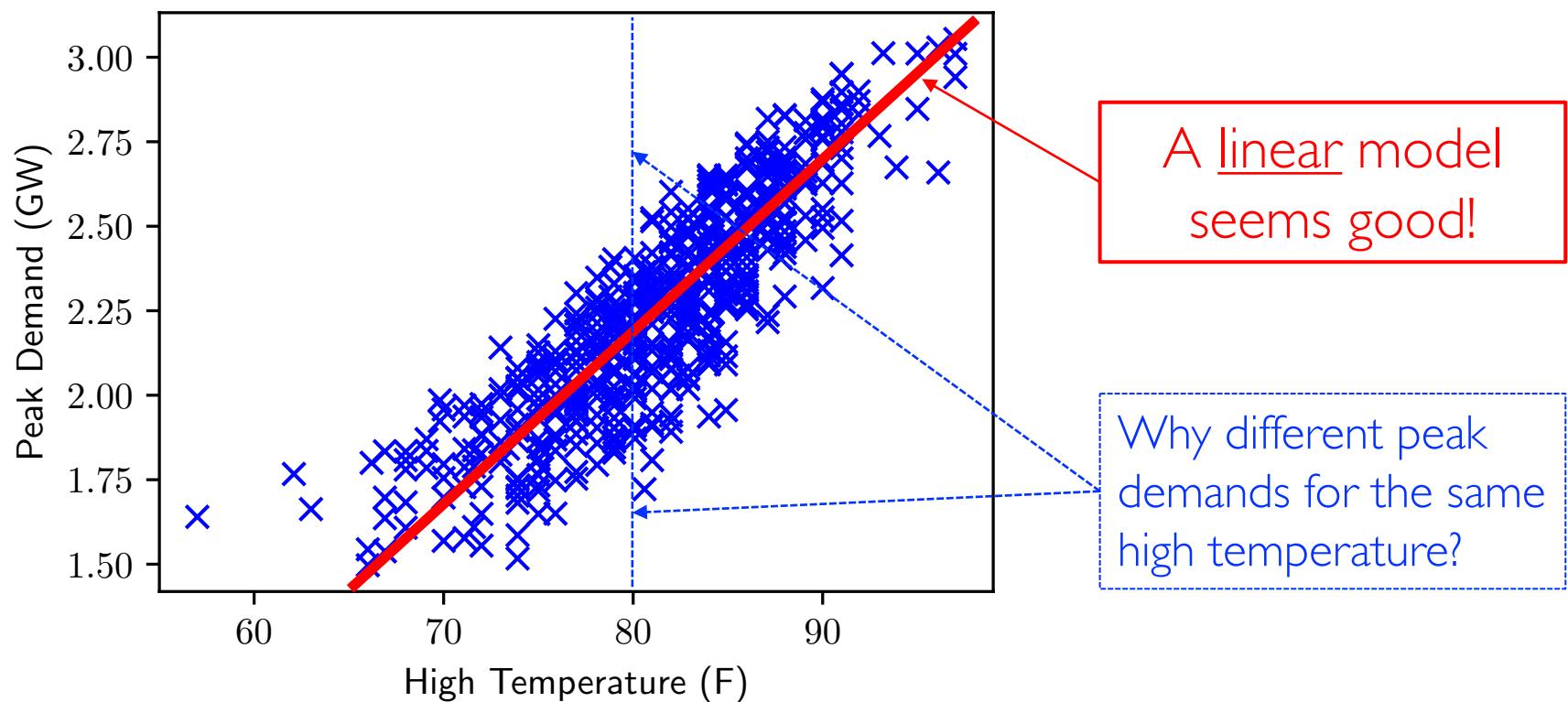
- Problem: Predict the peak power consumption in summer.



How to map from \mathcal{X} to \mathcal{Y} ?

What Model to Choose?

- Exploratory Data Analysis (EDA) to ease model selection.
 - Plot *high temperature* vs. *peak demand* in summer (June-August)



Linear Regression: Hypothesis Space

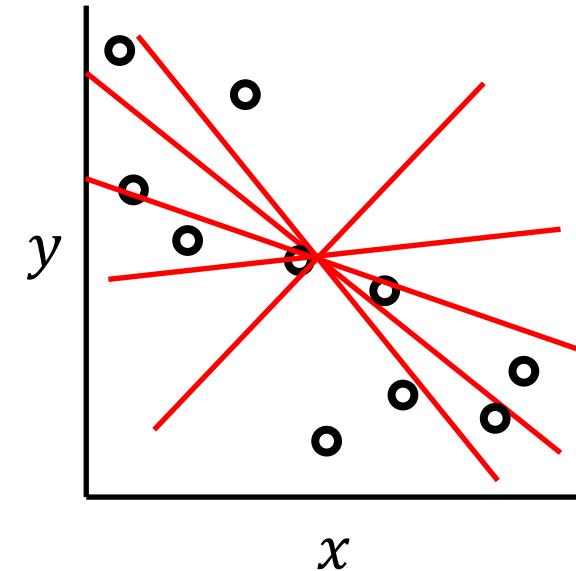
- Training data: $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n)$
 - Feature vector: $\mathbf{x} \in \mathbb{R}^d$, response: $y \in \mathbb{R}$
- Add a constant dimension to feature vector \mathbf{x} :

Equivalent to add intercept
in the linear models.

$$\mathbf{x} = \begin{pmatrix} 1 \\ x_1 \\ \dots \\ x_d \end{pmatrix} \in \mathbb{R}^{d+1}$$

- Prediction of hypothesis h parametrized by \mathbf{w} :

$$h(\mathbf{x}) = \boxed{w_0}_{\text{intercept}} + \sum_{j=1}^d w_j x_j = \sum_{j=0}^d w_j x_j = \boxed{\mathbf{w}} \cdot \mathbf{x}_{\text{normal vector}}$$



Each hypothesis
(high-dim plane)
corresponds to a \mathbf{w}

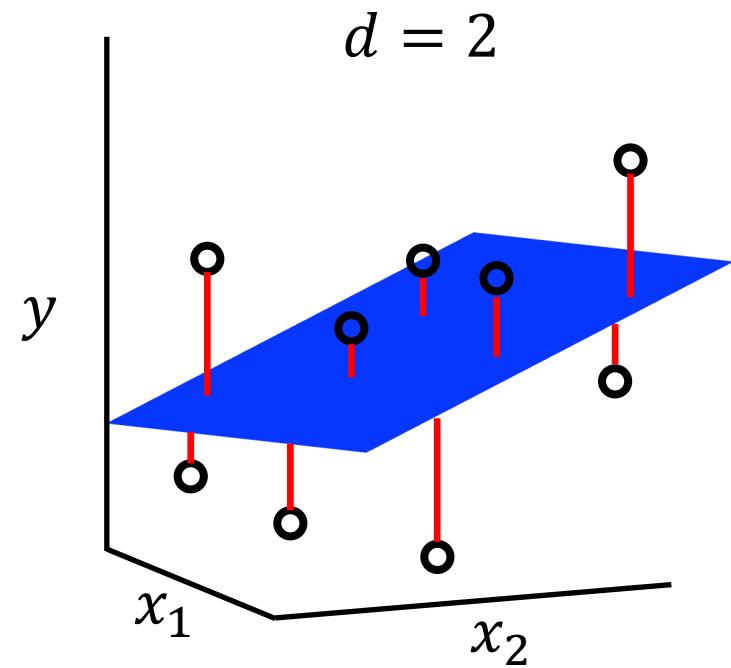
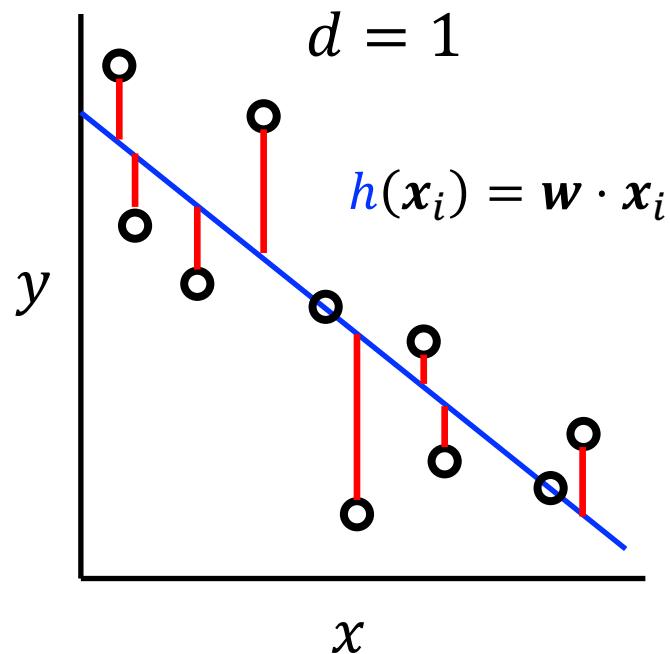
Linear Regression: Loss Function

- Use the squared loss (L2 loss) to compute the error on training set:

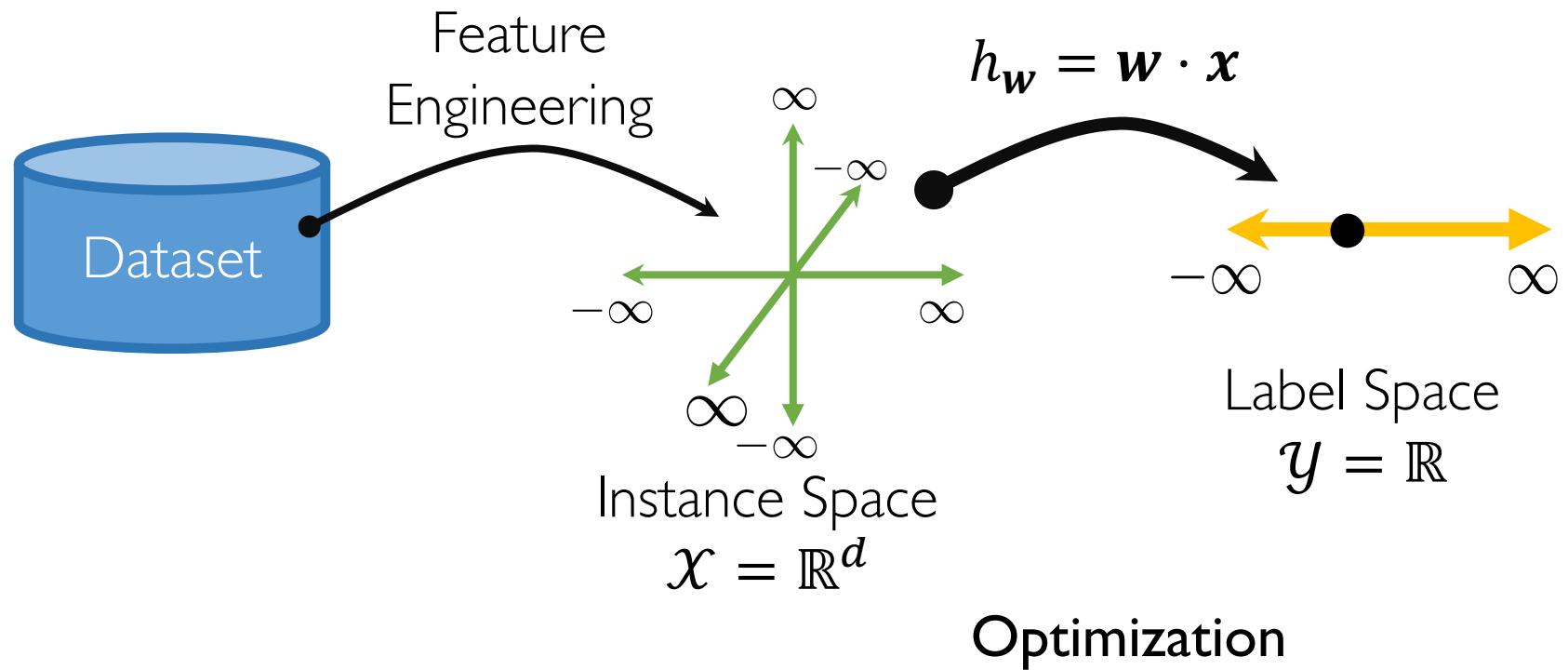
[^] standards for
empirical

$$\hat{\epsilon}(h) = \sum_{i=1}^n (h(x_i) - y_i)^2$$

Training Error



Linear Regression (线性回归)



How do we solve the
training problem
for the optimal \mathbf{w} ?

$$\min_{\mathbf{w}} \sum_{i=1}^n (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$$

Linear Regression: Analytical Solution

- Training error of linear regression:

$$\begin{aligned}\hat{\epsilon}(\mathbf{w}) &= \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 \\ &= \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2\end{aligned}$$

Matrix
Cookbook
is all you need!

- Put everything in matrix form:

data matrix
(design matrix)
 $n \times (d + 1)$

$$\mathbf{X} = \begin{bmatrix} -\mathbf{x}_1^T - \\ -\mathbf{x}_2^T - \\ \vdots \\ -\mathbf{x}_n^T - \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

label matrix
 $n \times 1$

<https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

Linear Regression: Analytical Solution

- Computing the gradient of $\hat{\epsilon}(\mathbf{w})$ w.r.t. \mathbf{w} and setting it to zero yields the optimal parameter \mathbf{w}^* :

$$\hat{\epsilon}(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

In Matrix Cookbook:

- $\|\mathbf{A}\|_F^2 = \text{tr}(\mathbf{A}^T \mathbf{A})$
- $\frac{\partial \text{tr}(\mathbf{A}^T \mathbf{B})}{\partial \mathbf{A}} = \mathbf{B}$

$$\begin{aligned}\nabla_{\mathbf{w}} \hat{\epsilon}(\mathbf{w}) &= 2\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{0} \\ &\Rightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \\ &\Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

Normal Equation

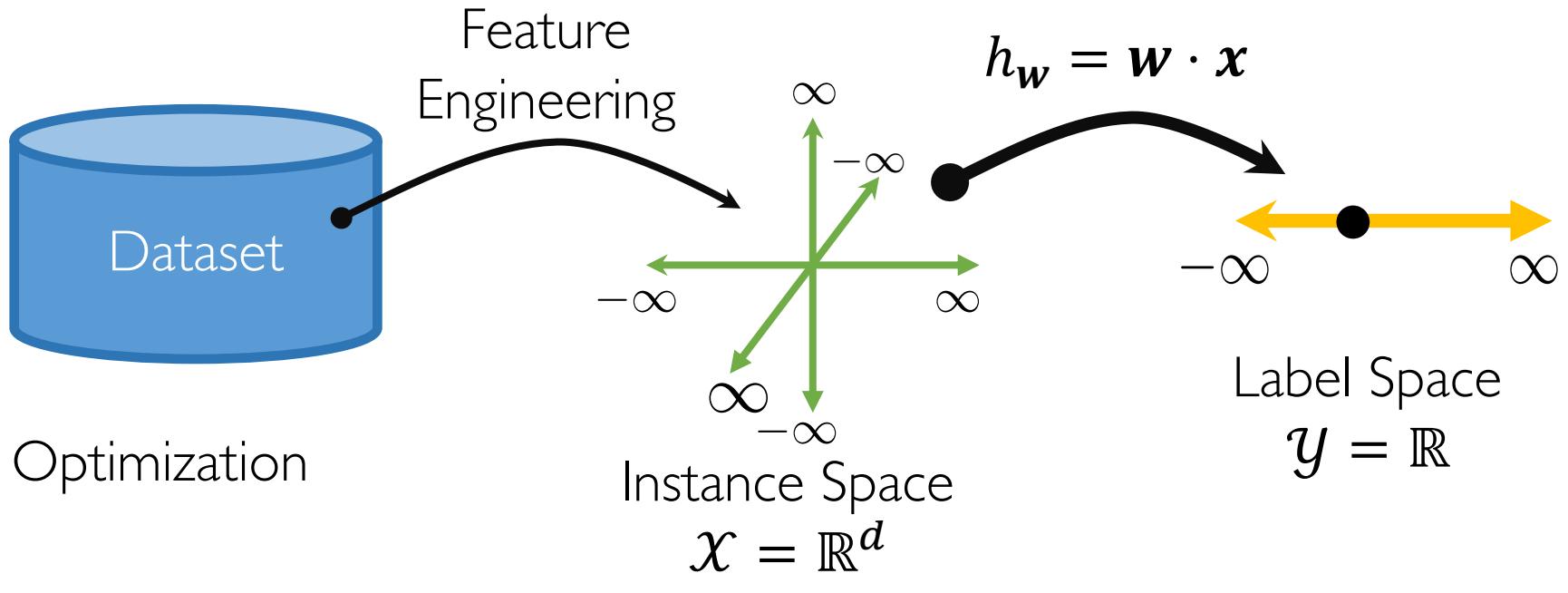
$$\left(\underbrace{\begin{pmatrix} \mathbf{X}^T \mathbf{X} \end{pmatrix}}_{(d+1) \times (d+1)} \right)^{-1} \times \underbrace{\begin{pmatrix} \mathbf{X}^T \end{pmatrix}}_{(d+1) \times n}$$

$(d+1) \times n$

Computational complexity is **big!**

$O[d^2(d+n)]$

Linear Regression



Optimization

Analytic Solution

$$w = (X^T X)^{-1} X^T y$$

Unaffordable cost for
high-dimensional big data!

$$\min_w \sum_{i=1}^n (h_w(x_i) - y_i)^2$$

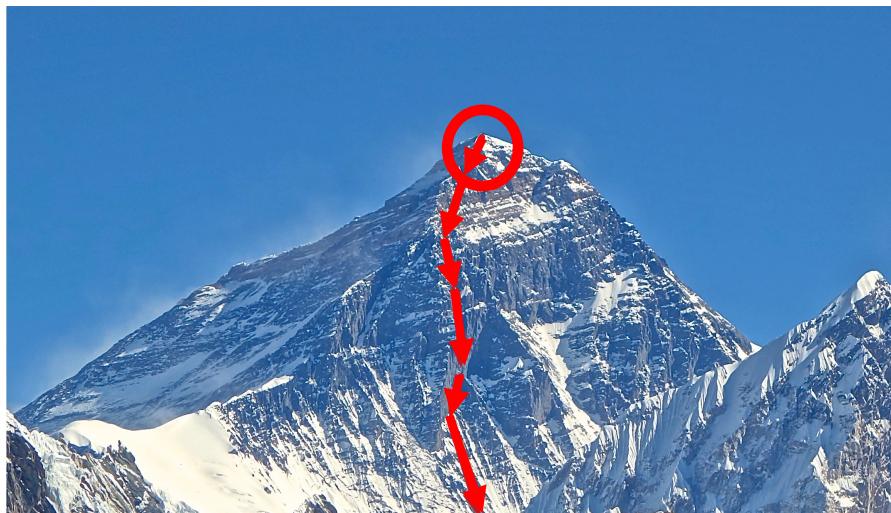
What should we do in practice?

Outline

- Linear Regression
 - Optimization
 - Statistical View: Maximum Likelihood Estimation
- Nonlinearization
- Regularization
 - L1 and L2 Norm
 - Statistical View: Maximum a Posteriori Estimation
- Linear Classification
 - Logistic Regression
 - Softmax Regression

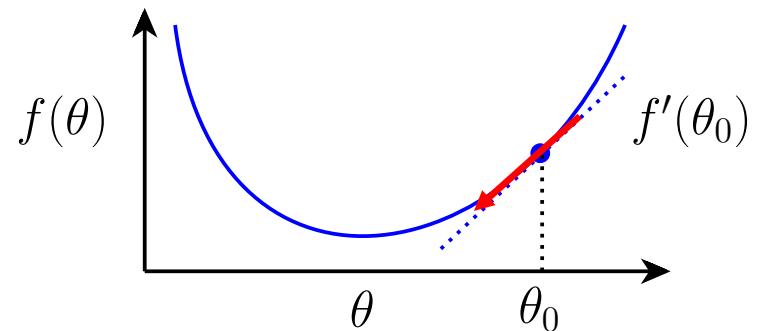
Optimization

- Normal equation suffers from **high computing cost**.
- Loss functions for other tasks are **hard to derive an analytic solution**.
- For general differentiable loss function, use **Gradient Descent (GD)**



Follow the steepest slope

You only need to observe local information for each step.



Work well for convex function

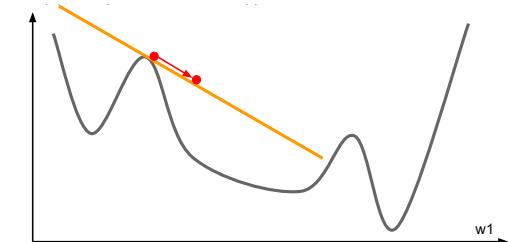
Gradient (梯度)

- Gradient shows direction that function varies fastest.

Same dimension with
parameter \mathbf{w} 's dimension

$$\mathbf{g} = \nabla_{\mathbf{w}} J(\mathbf{w})$$

$$g_j = \nabla_{w_j} J(\mathbf{w})$$



- First-order Taylor approximation of the objective function:

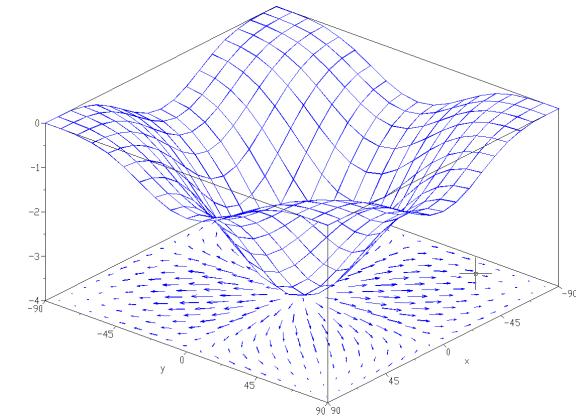
$$J(\mathbf{w}) = J(\mathbf{w}_0) + (\mathbf{w} - \mathbf{w}_0)^T \mathbf{g} + \dots$$

- Go along gradient for a step with a small rate η :

$$J(\mathbf{w} - \eta \mathbf{g}) \approx J(\mathbf{w}) - \eta \mathbf{g}^T \mathbf{g}$$

- Reach a point with smaller loss.

$$-\eta \mathbf{g}^T \mathbf{g} \leq 0$$



- Repeat this step, and we get the Gradient Descent (GD) algorithm.

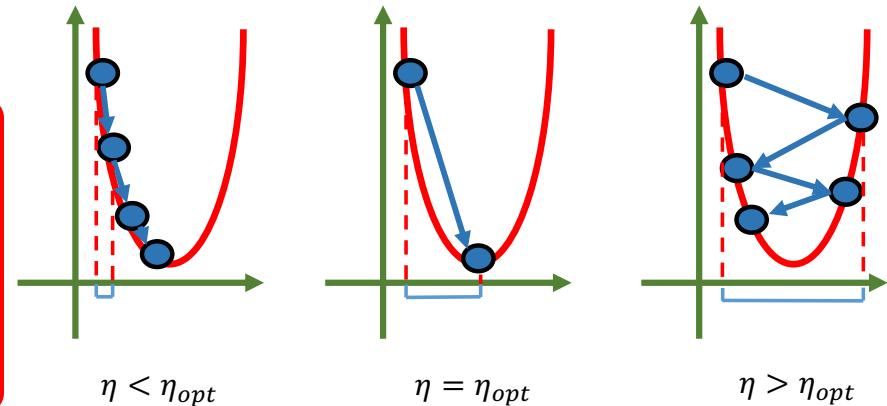
Gradient Descent (GD, 梯度下降)

- General optimization algorithm

- For each iteration t ($\leq T$)

- When \mathbf{w}^t does not converge:
- $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta \nabla_{\mathbf{w}} \hat{\epsilon}(\mathbf{w})|_{\mathbf{w}=\mathbf{w}^t}$

- Learning rate η matters!



Squared Loss is convex.

- For linear regression:

- $\nabla_{\mathbf{w}} \hat{\epsilon}(\mathbf{w}) = 2\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y})$

- When \mathbf{w}^t does not converge:

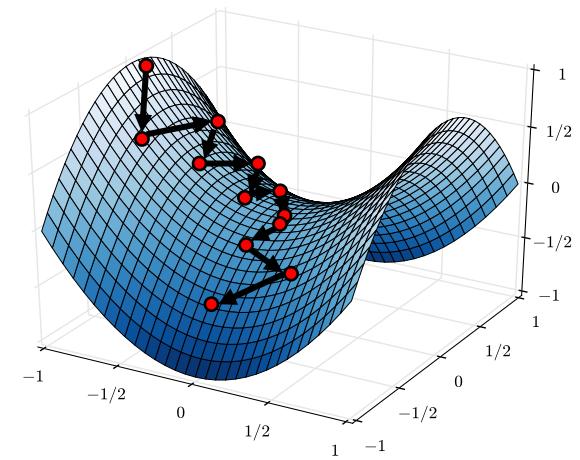
- $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - 2\eta \mathbf{X}^T(\mathbf{X}\mathbf{w}^t - \mathbf{y})$

No need for computing inverse matrix!

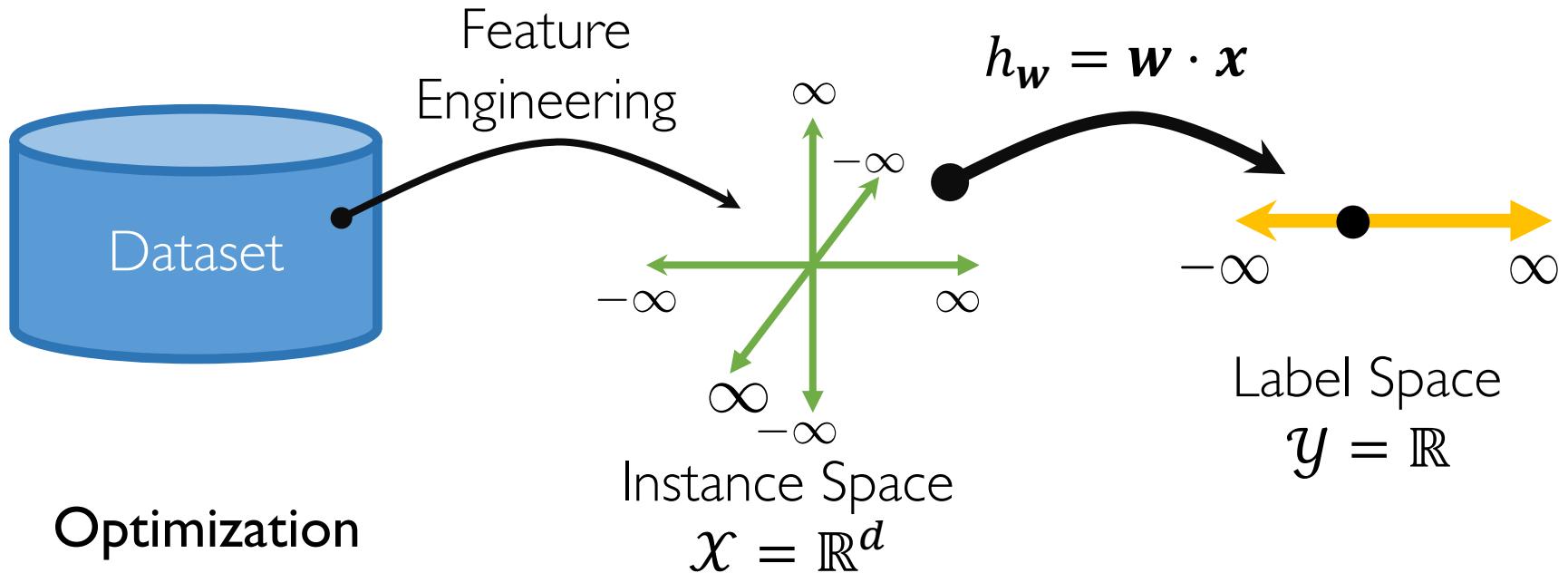
Complexity (for T iterations):
 $O(dnT)$

Stochastic Gradient Descent (SGD)

- In each iteration t ($\leq T$):
 - Randomly sample a minibatch of $m \ll n$ points $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$
 - Set $J^t(\mathbf{w}^t) = \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{w}^t; \mathbf{x}_i, y_i)$
 - Compute gradient on minibatch:
$$\Delta^t = \nabla_{\mathbf{w}} J^t(\mathbf{w}^t)$$
 - Update parameters with learning rate η :
$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \Delta^t$$
- For linear regression: $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - 2\eta \mathbf{X}_m^T (\mathbf{X}_m \mathbf{w}^t - \mathbf{y})$
 - Computational complexity $O(dmT)$. What if large dimension d ?



Linear Regression



Why should we use squared loss?

Outline

- Linear Regression
 - Optimization
 - Statistical View: Maximum Likelihood Estimation
- Nonlinearization
- Regularization
 - L1 and L2 Norm
 - Statistical View: Maximum a Posteriori Estimation
- Linear Classification
 - Logistic Regression
 - Softmax Regression

Likelihood in Parametric Models

- Suppose we have a parametric model $\{p(\mathbf{z}; \theta) | \theta \in \Theta\}$ parametrized by θ and an i.i.d. sample dataset $\mathcal{D} = (\mathbf{z}_1, \dots, \mathbf{z}_n)$
- The likelihood of the estimated parameter $\hat{\theta} \in \Theta$ for sample \mathcal{D} is

$$p(\mathcal{D}; \hat{\theta}) \triangleq \prod_{i=1}^n p(\mathbf{z}_i; \hat{\theta})$$

- Due to numerical stability, we prefer to work with the log-likelihood:

$$\log p(\mathcal{D}; \hat{\theta}) \triangleq \sum_{i=1}^n \log p(\mathbf{z}_i; \hat{\theta})$$

- And sums are easier to work with than products.

Maximum Likelihood Estimation (MLE)

- Maximum Likelihood Estimator (MLE) for estimating the parameter θ in the parametric model $\{p(\mathbf{z}; \theta) | \theta \in \Theta\}$ is

$$\hat{\theta} \in \operatorname{argmax}_{\theta \in \Theta} \log p(\mathcal{D}; \theta) \Rightarrow \hat{\theta} \in \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log p(\mathbf{z}_i; \theta)$$

- Finding the MLE is an optimization problem.

- Under some proper distributional assumption:

- The MLE of a parametric model leads to a particular loss function.
- Note that: We need to specify the distribution before solving MLE.

Think of this for all models.

Linear Regression: Statistical View

- In the **supervised learning** model, $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ are **i.i.d.** samples from some data-generating distribution $p(\mathbf{x}, y)$.

$$p(\mathcal{D}; \hat{\theta}) \triangleq \prod_{i=1}^n p(\mathbf{x}_i, y_i; \hat{\theta}) = \prod_{i=1}^n p(y_i | \mathbf{x}_i; \hat{\theta}) p(\mathbf{x}_i)$$

- For the **discriminative model** that only models the relation between \mathbf{x} and y , we assume \mathbf{x}_i follows an **empirical** distribution $p(\mathbf{x}_i) = \frac{1}{n}$.

$$\prod_{i=1}^n p(y_i | \mathbf{x}_i; \hat{\theta}) p(\mathbf{x}_i) = \prod_{i=1}^n \frac{1}{n} p(y_i | \mathbf{x}_i; \hat{\theta}) \propto \prod_{i=1}^n p(y_i | \mathbf{x}_i; \hat{\theta})$$

Linear Regression: Statistical View

- In the **discriminative model**, the likelihood and the log-likelihood of the estimated parameter $\hat{\theta} \in \Theta$ for sample \mathcal{D} :

$$p(\mathcal{D}; \theta) \triangleq \prod_{i=1}^n p(y_i | \mathbf{x}_i; \theta)$$

$$\log p(\mathcal{D}; \theta) = \sum_{i=1}^n \log p(y_i | \mathbf{x}_i; \theta)$$

- The Maximum Likelihood Estimation (MLE) for **discriminative model**:

$$\hat{\theta} = \max_{\theta} \sum_{i=1}^n \log p(y_i | \mathbf{x}_i; \theta)$$

Linear Regression: Statistical View

- In linear regression problem, we assume that:

$$y \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}, \sigma^2)$$

$$y = \mathbf{w}^T \mathbf{x} + e$$

Gaussian noise assumption

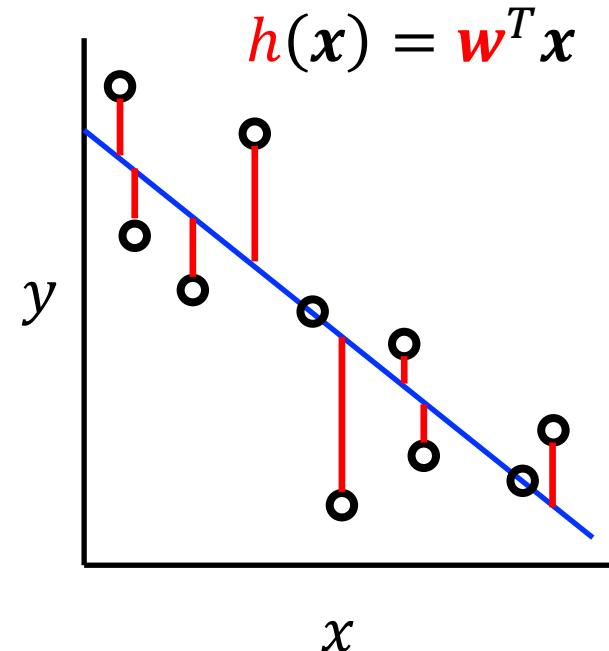
$$e \sim \mathcal{N}(\mathbf{0}, \sigma^2)$$

where y 's are independent.

- The linear regression should correctly model:

$$\mathbb{E}(y|\mathbf{x}; \mathbf{w}, \sigma^2) = \mathbf{w}^T \mathbf{x}$$

Regression function



- We have i.i.d. sampled data $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$.
- We can find the best parameter \mathbf{w}^* for sample \mathcal{D} through the MLE.

Linear Regression: Statistical View

- Assume that $y_i|\mathbf{w}, \mathbf{x}_i \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}_i, \sigma^2)$, then $\theta = \{\mathbf{w}, \sigma\}$.

- For each point (\mathbf{x}_i, y_i) :

$$p(y_i|\mathbf{x}_i; \mathbf{w}, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \mathbf{w}^T \mathbf{x}_i)^2\right\}$$

- The log-likelihood for linear regression on the whole dataset \mathcal{D}_n :

$$\begin{aligned} \log p(\mathcal{D}_n; \mathbf{w}, \sigma) &= \sum_{i=1}^n \log p(y_i|\mathbf{x}_i; \mathbf{w}, \sigma) \\ &= \frac{n}{2} \log \frac{1}{2\pi\sigma^2} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \end{aligned}$$

Linear Regression: Statistical View

- The Maximum Likelihood Estimation of the optimal parameter \mathbf{w}^* :

$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} \frac{n}{2} \log \frac{1}{2\pi\sigma^2} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

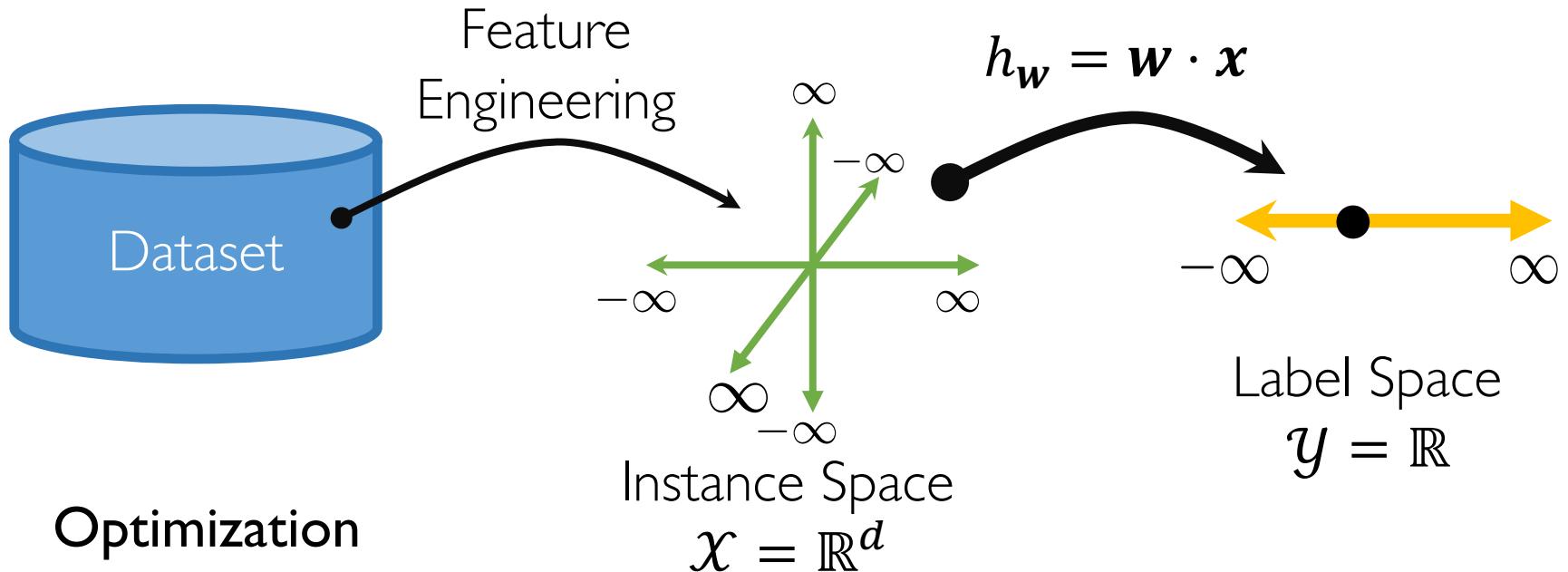
$$= \operatorname{argmax}_{\mathbf{w}} -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

$$= \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

How about estimating σ ?
How about other noise model?

- Note: Maximum Likelihood Estimation of \mathbf{w} under the Gaussian noise assumption is equivalent to linear regression with the **squared loss**

Linear Regression



Optimization

Gradient Descent

$$\nabla_{\mathbf{w}} \hat{\epsilon}(\mathbf{w}) = 2\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y})$$

Analytic Solution

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

$$\min_{\mathbf{w}} \sum_{i=1}^n (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$$

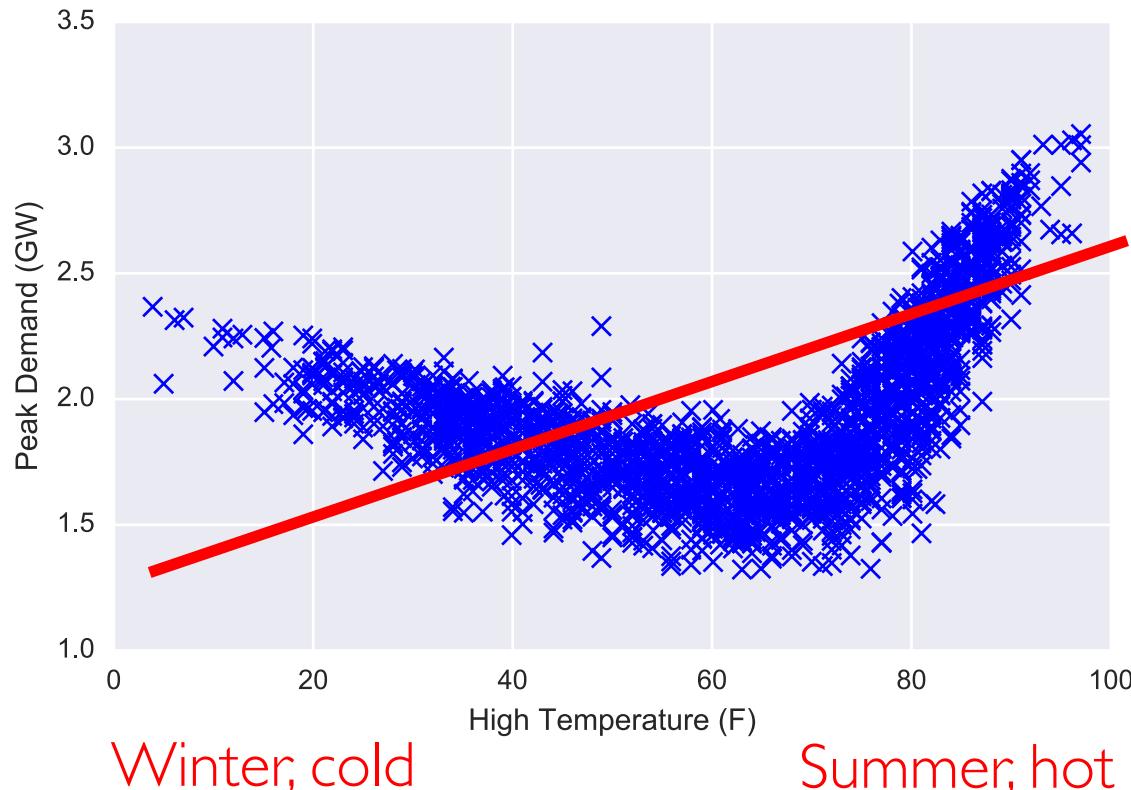
Loss Function
How to handle nonlinearity?

Outline

- Linear Regression
 - Optimization
 - Statistical View: Maximum Likelihood Estimation
- **Nonlinearization**
- Regularization
 - L1 and L2 Norm
 - Statistical View: Maximum a Posteriori Estimation
- Linear Classification
 - Logistic Regression
 - Softmax Regression

An Example: Electricity Prediction

- Problem: Predict the peak power consumption in all months.
- Exploratory Data Analysis (EDA): Peak demand vs. temperature plot



Can we use linear regression again?

Underfit!

Basis Function (基函数)

- Feature map: $\mathbf{x} = (x_1, \dots, x_d) \in \mathcal{X} \xrightarrow{\Phi} \mathbf{z} = (z_1, \dots, z_{\tilde{d}}) \in \mathcal{Z}$

$$\mathbf{z} = \Phi(\mathbf{x})$$

- Each $z_j = \phi_j(\mathbf{x})$ depends on some nonlinear transform ϕ_j .

- $\{\phi_j\}_{0 \leq j \leq \tilde{d}}$ is called **basis functions**.

- Polynomial basis functions

1-D vector:

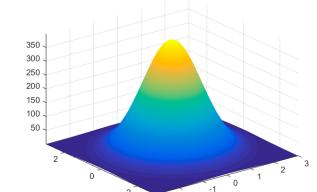
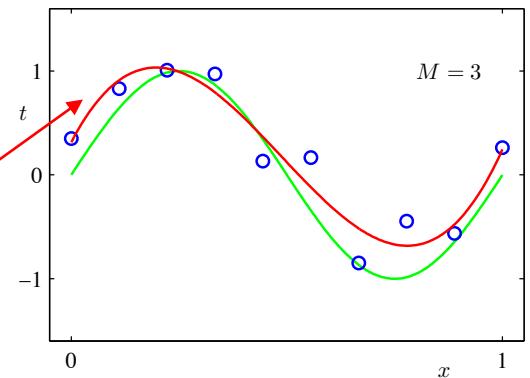
$$\mathbf{z}' = (1, x_1, x_1^2, x_1^3)$$

2-D vector:

$$\mathbf{z}' = (1, x_2, x_1 x_2, x_1^2)$$

- Radial basis functions (RBF)

$$\phi_j(\mathbf{x}) = \left\{ \exp \left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}_j\|_2^2}{2\sigma^2} \right) : j = 1, \dots, \tilde{d} \right\}$$

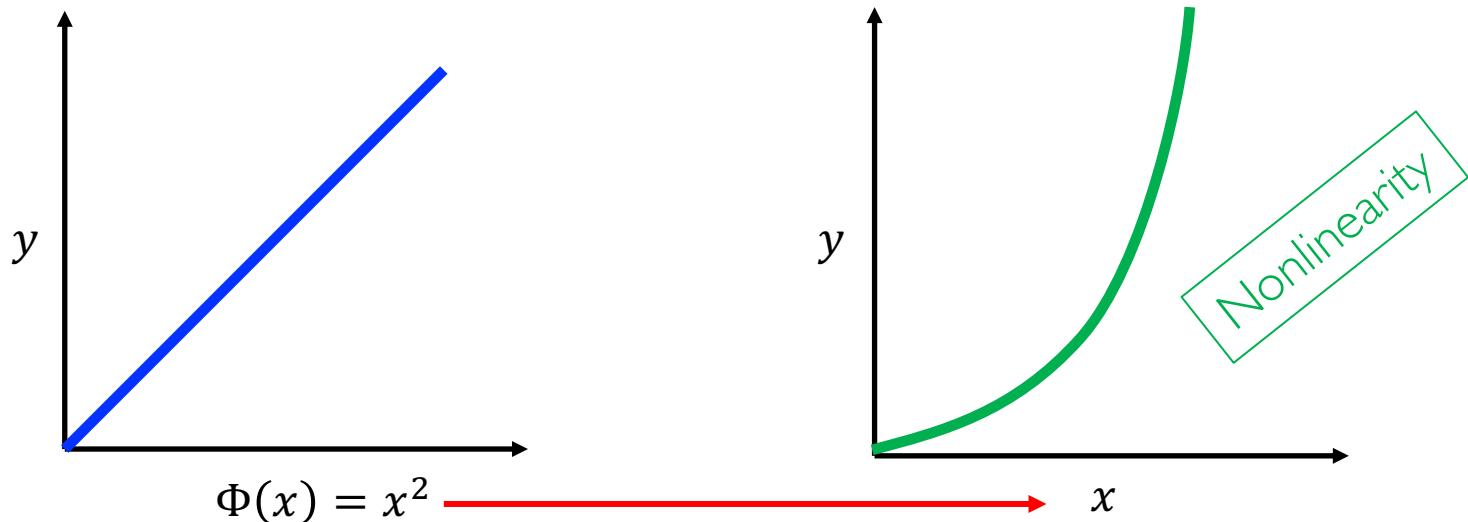


Linear Regression: Basis Function

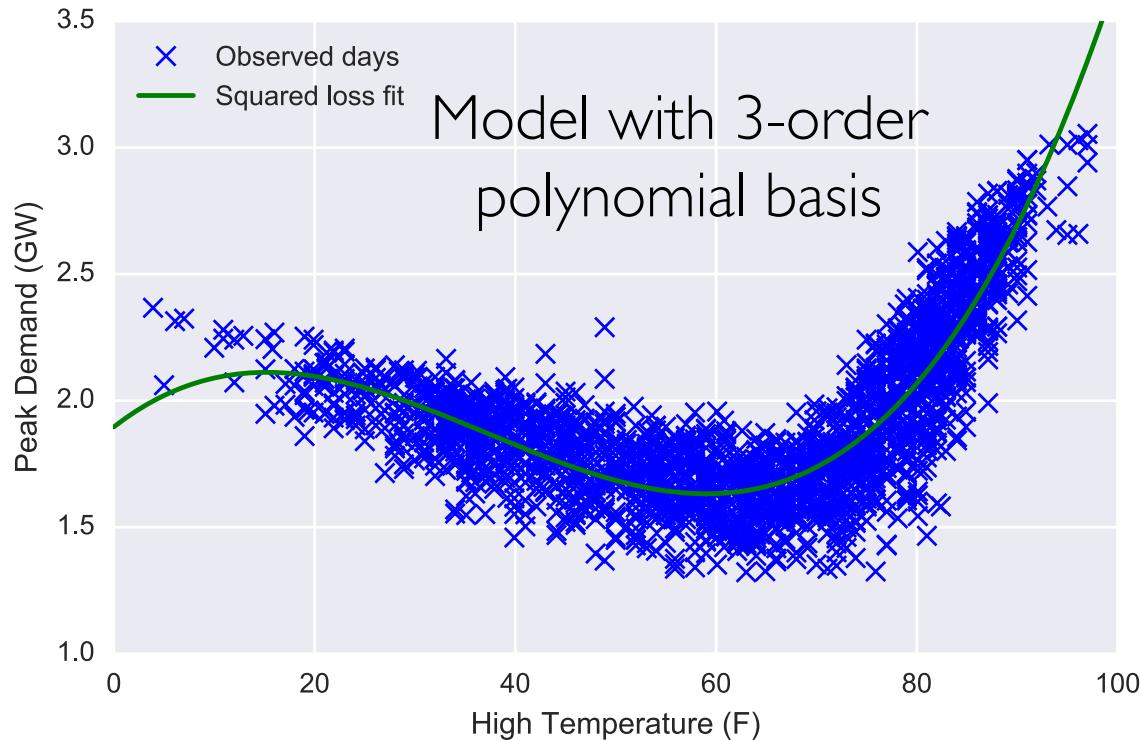
- The final hypothesis is linear in the feature space \mathcal{Z} .
- The final hypothesis is nonlinear in the input space \mathcal{X} :

$$h(\mathbf{x}) = \tilde{\mathbf{w}} \cdot \mathbf{z} = \tilde{\mathbf{w}} \cdot \Phi(\mathbf{x}) = \sum_{j=1}^{\tilde{d}} \tilde{w}_j \phi_j(\mathbf{x})$$

- Algorithm works because of linearity in the parameters (still simple!)

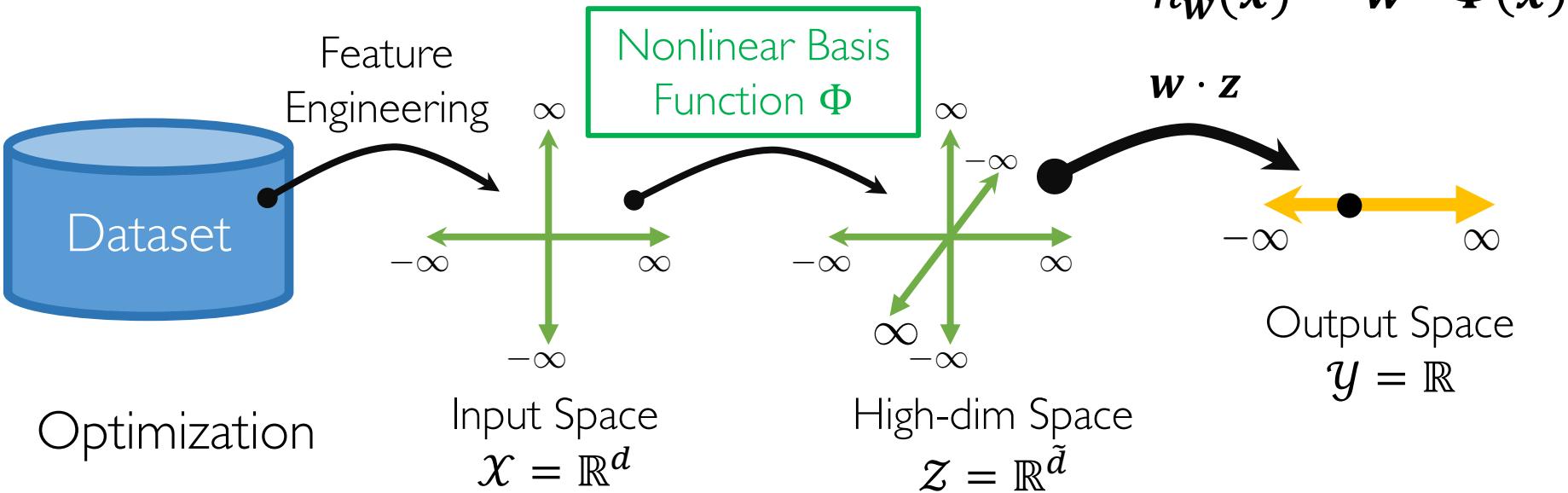


Linear Regression: Basis Function



- With polynomial basis functions, we could **fit** nonlinear dataset.
 - Linear regression with polynomial basis functions.
- Higher-order polynomial has **stronger** ability to fit complex data.

Linear Regression



Gradient Descent

$$\nabla_{\mathbf{w}} \hat{\epsilon}(\mathbf{w}) = 2\mathbf{Z}^T(\mathbf{Z}\mathbf{w} - \mathbf{y})$$

Analytic Solution

$$\mathbf{w} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y}$$

Loss Function

$$\min_{\mathbf{w}} \sum_{i=1}^n (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$$

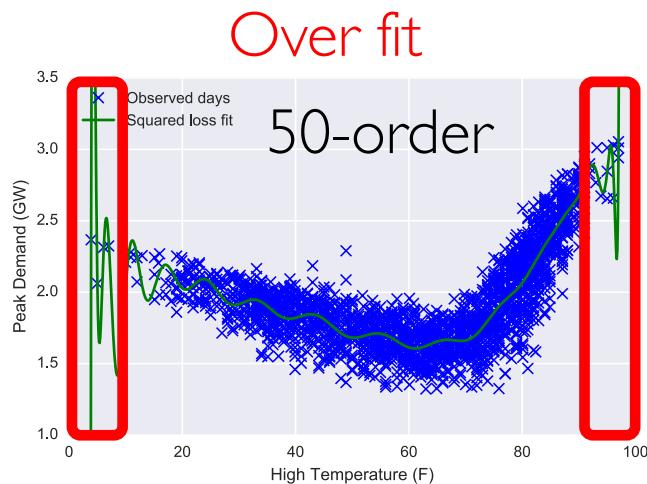
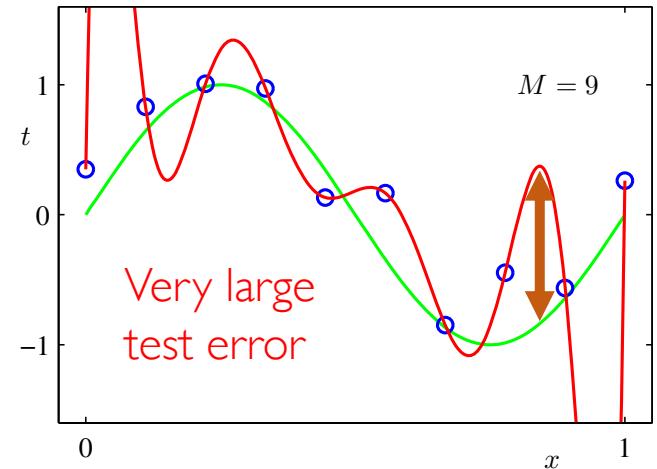
How to deal with overfitting?

Outline

- Linear Regression
 - Optimization
 - Statistical View: Maximum Likelihood Estimation
- Nonlinearization
- **Regularization**
 - L1 and L2 Norm
 - Statistical View: Maximum a Posteriori Estimation
- Linear Classification
 - Logistic Regression
 - Softmax Regression

Overfitting

- Higher-order polynomial fits data better.
- But may **fit noise!**
- Also may cause **very large coefficients**.



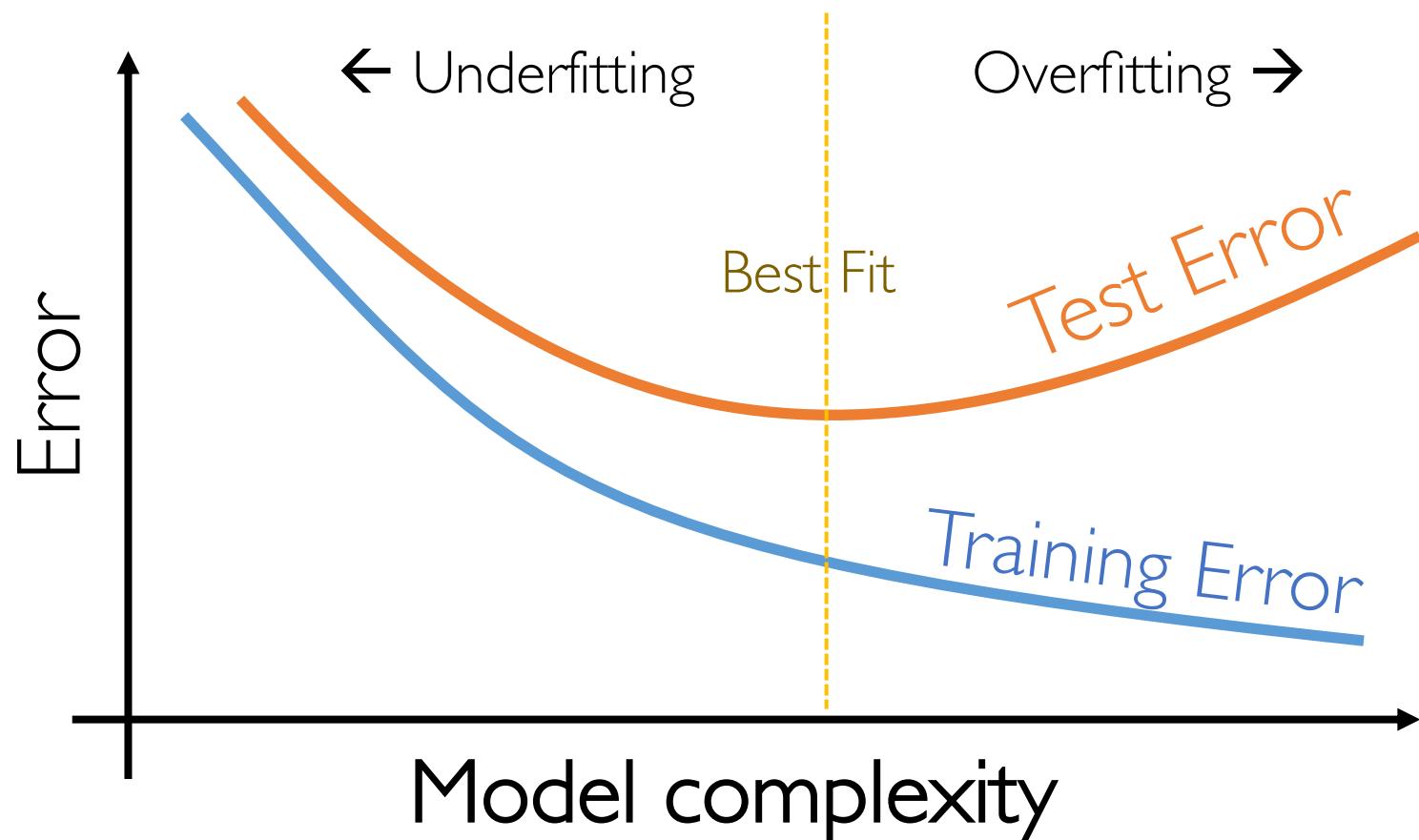
w_0^*	0.35
w_1^*	232.37
w_2^*	-5321.83
w_3^*	48568.31
w_4^*	-231639.30
w_5^*	640042.26
w_6^*	-1061800.52
w_7^*	1042400.18
w_8^*	-557682.99
w_9^*	125201.43

Solution: Control norm of $\mathbf{w}!$

Why?

When the size of dataset is fixed:

Training vs Test Error



How to control the model complexity? Solution: Control norm of \mathbf{w} !

Linear Regression: L2-Regularization

- The L2-regularized linear regression (ridge regression) imposes L2-norm regularization trading off with a hyperparameter $\lambda \geq 0$:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2$$

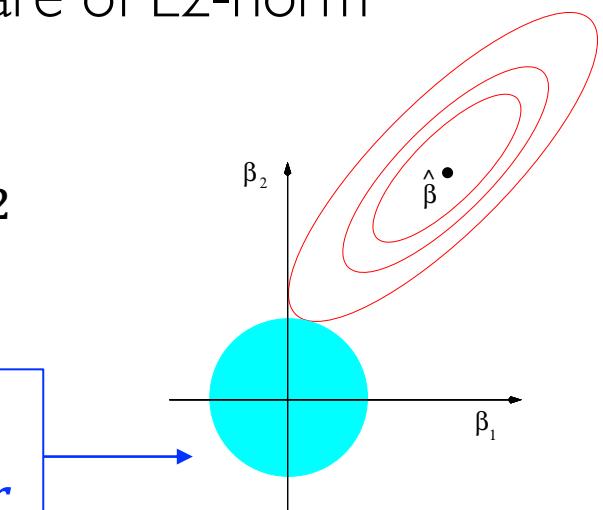
- Where $\|\mathbf{w}\|_2^2 = w_1^2 + \dots + w_d^2$ is the square of L2-norm
- Which is equivalent to:

$$\hat{\mathbf{w}} = \arg \min_{\|\mathbf{w}\|_2^2 \leq r} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$

for some $r > 0$.

- $r \rightarrow 0 \Leftrightarrow \lambda \rightarrow +\infty$

L2 unit ball
with radius r



Linear Regression: L1-Regularization

- The L1-regularized linear regression (lasso regression) imposes L1-norm regularization trading off with a hyperparameter $\lambda \geq 0$:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_1$$

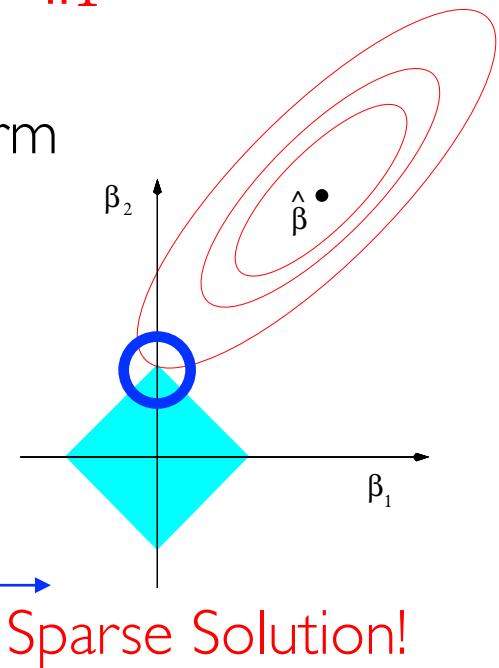
- Where $\|\mathbf{w}\|_1 = |w_1| + \dots + |w_d|$ is the L1-norm
- Which is equivalent to:

$$\hat{\mathbf{w}} = \arg \min_{\|\mathbf{w}\|_1 \leq r} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$

with some $r > 0$.

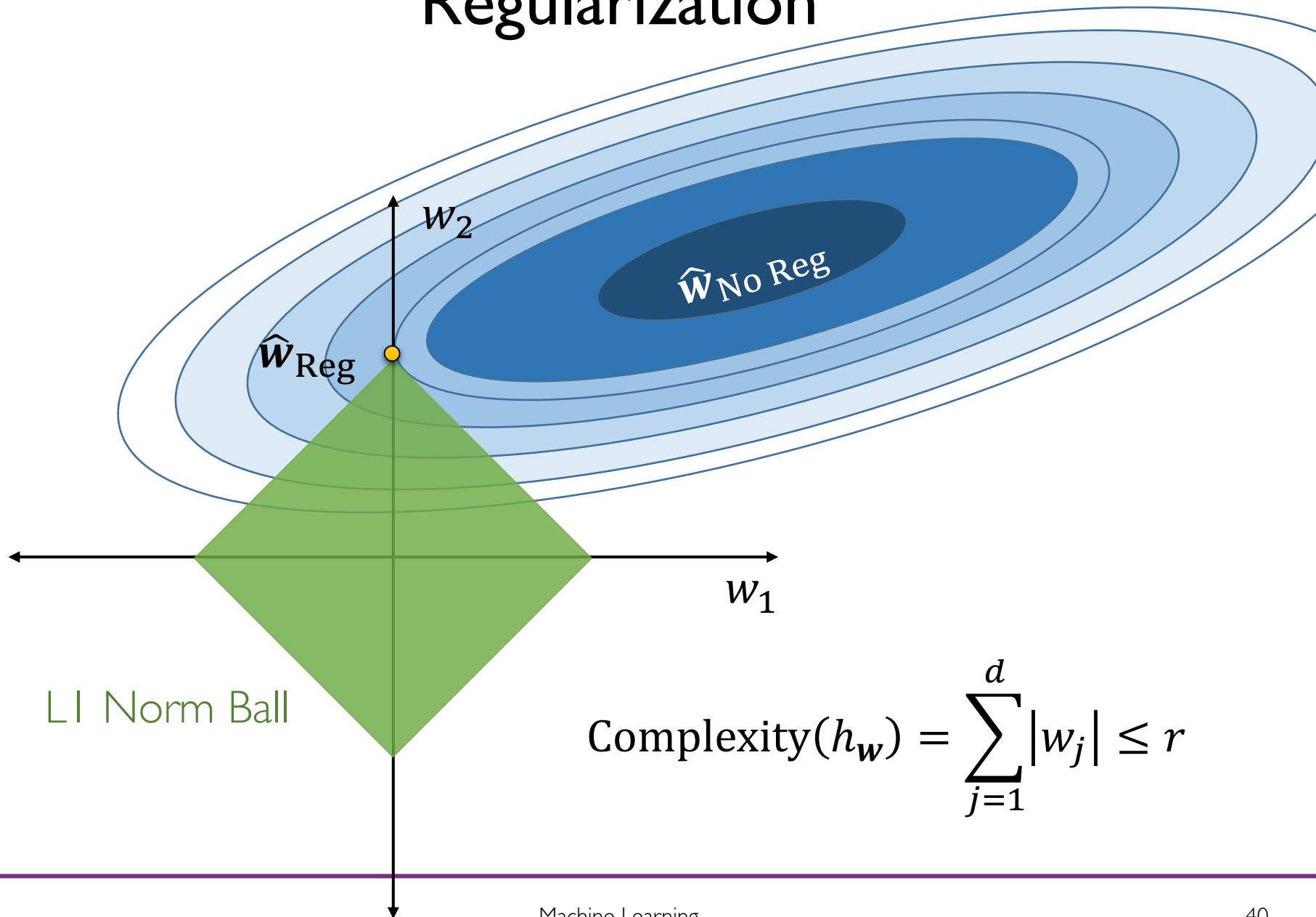
- $r \rightarrow 0 \Leftrightarrow \lambda \rightarrow +\infty$

L1 unit ball
with radius r

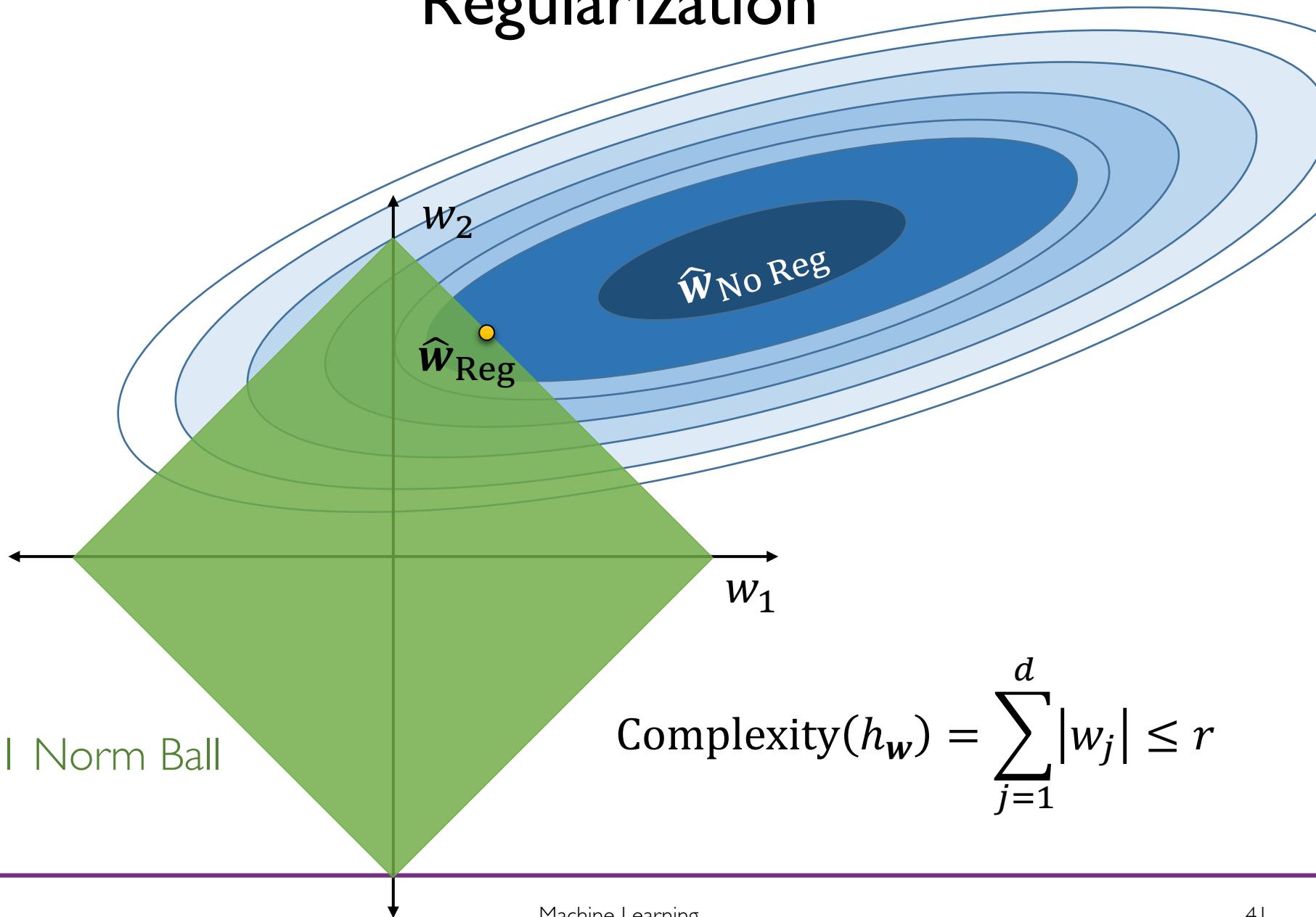


Sparse Solution!

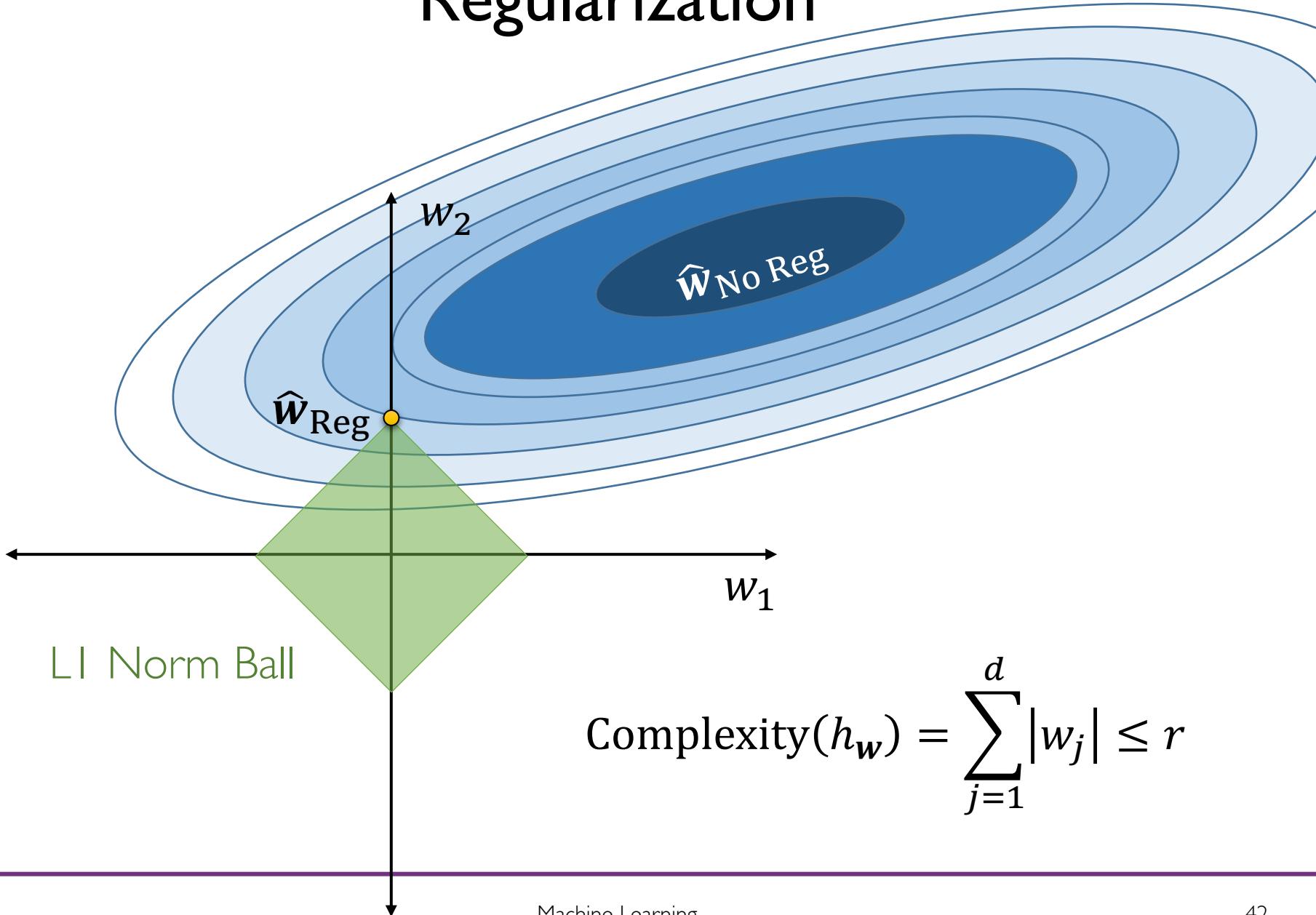
Regularization



Regularization

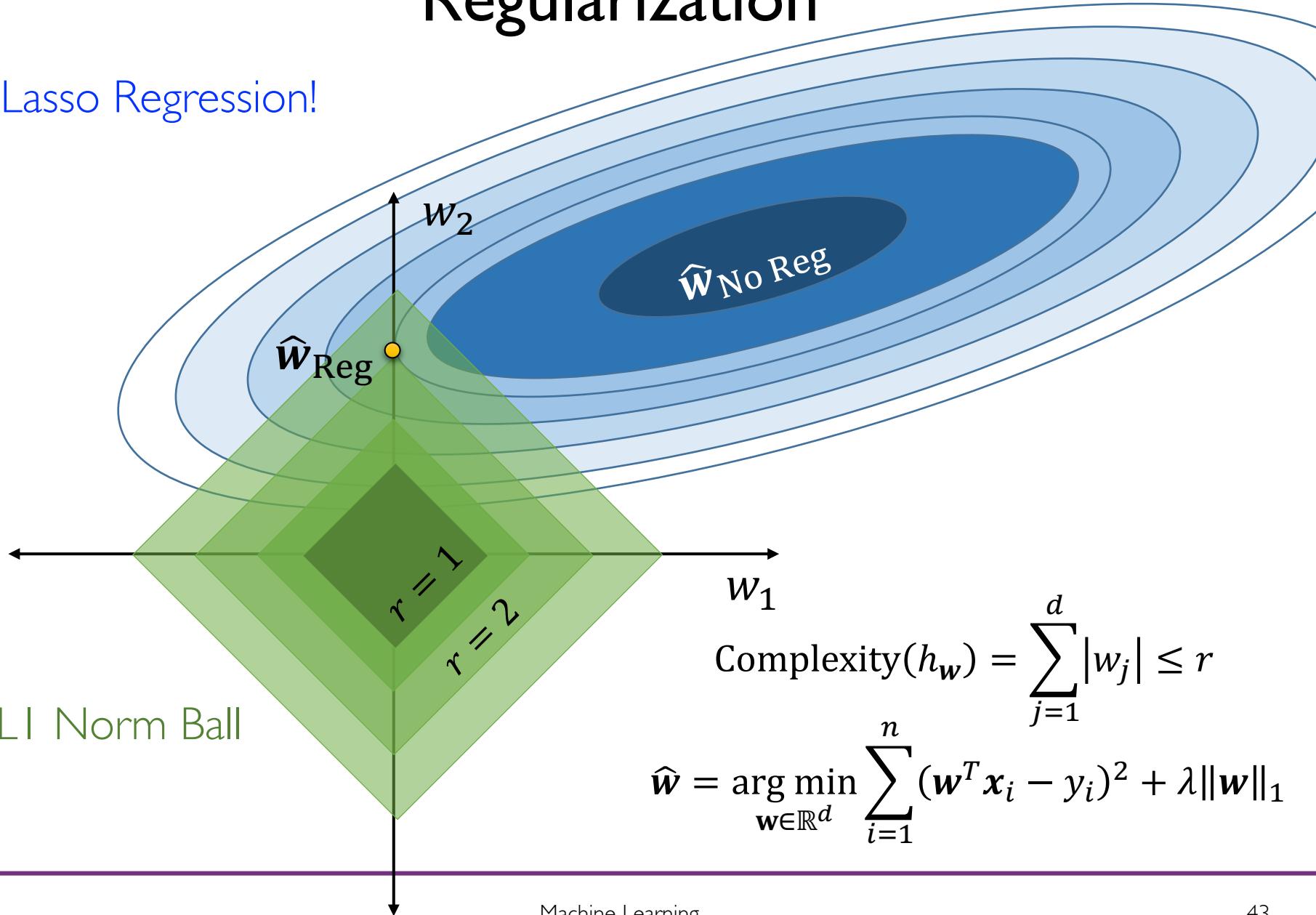


Regularization

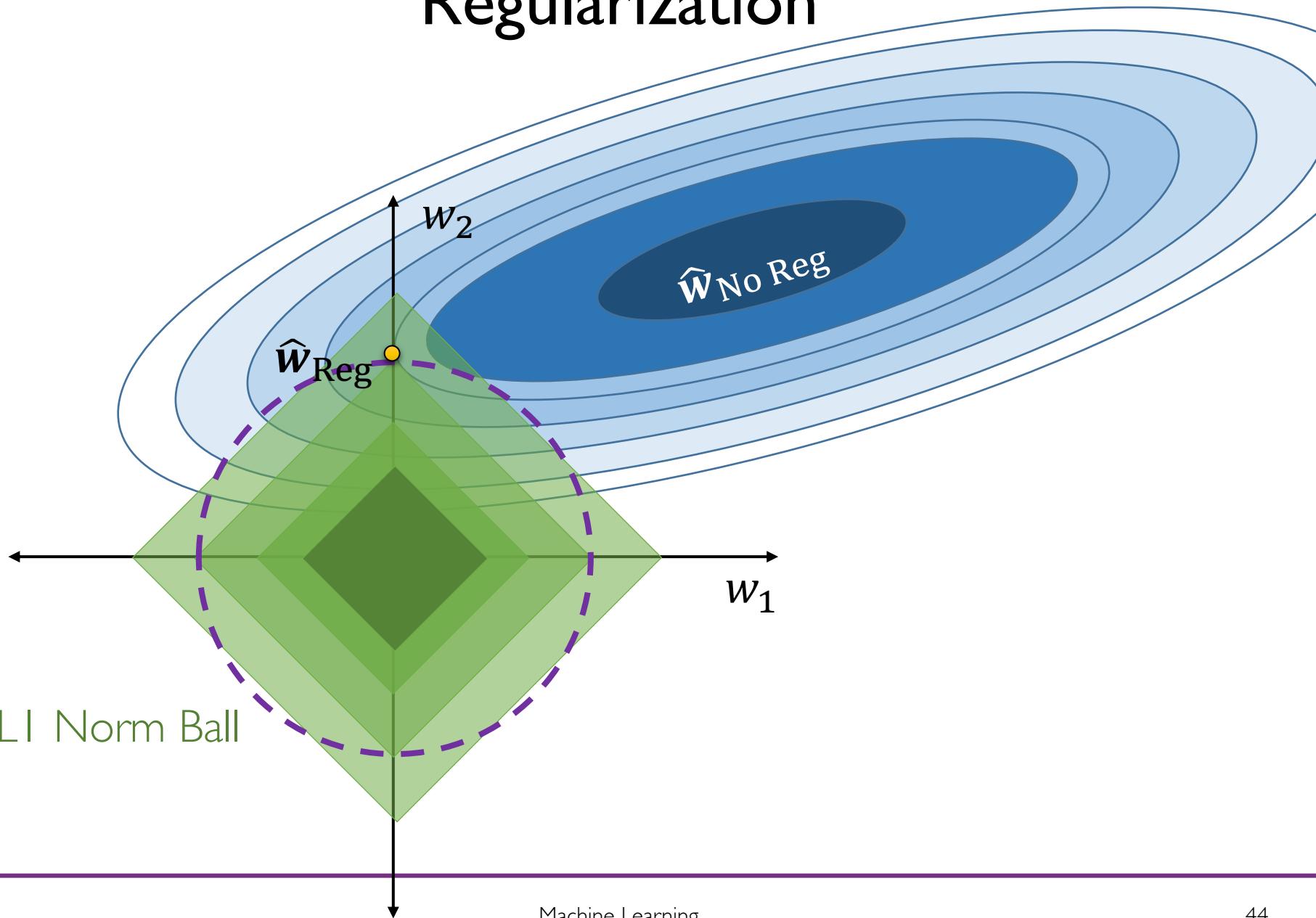


Regularization

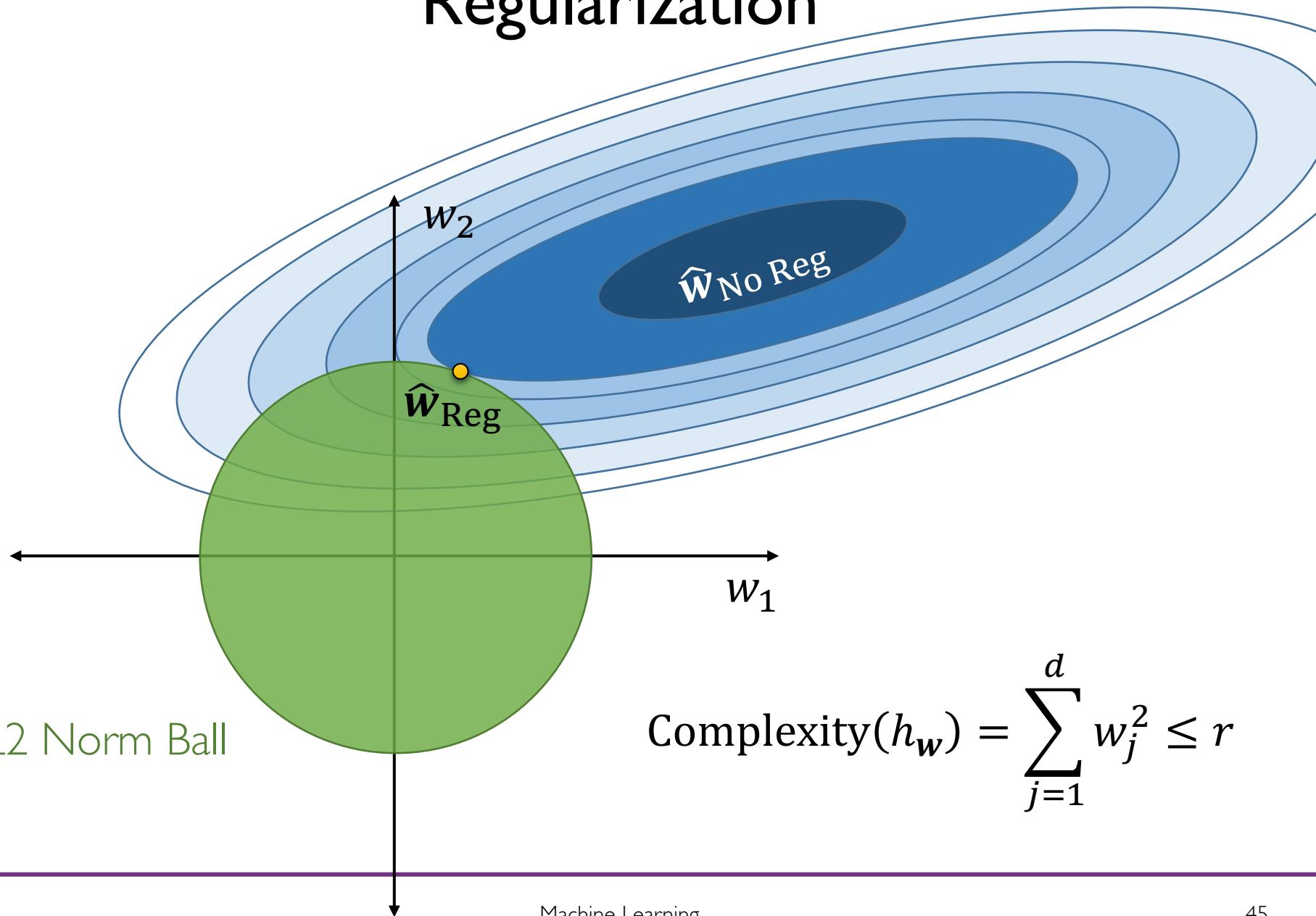
Lasso Regression!



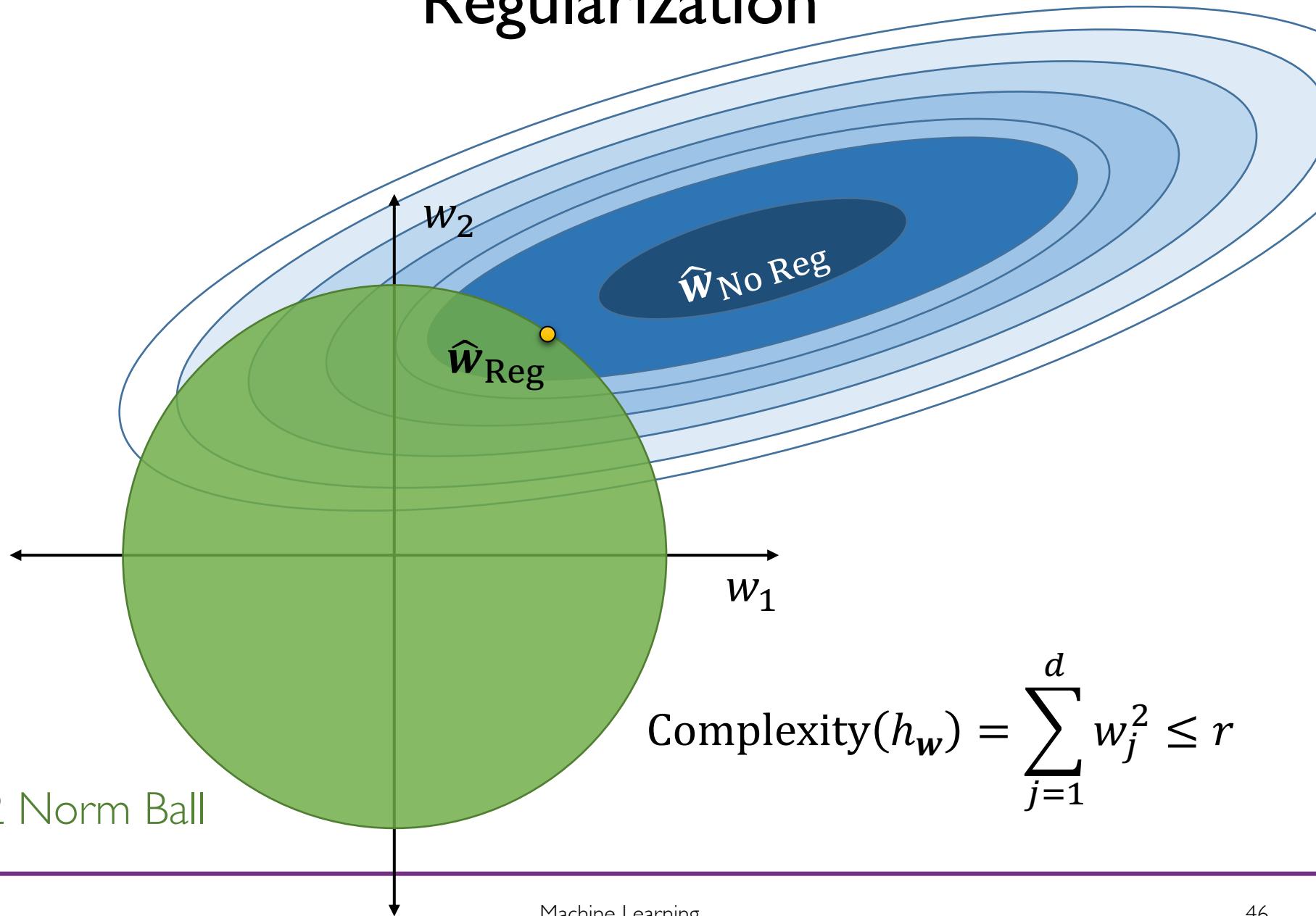
Regularization



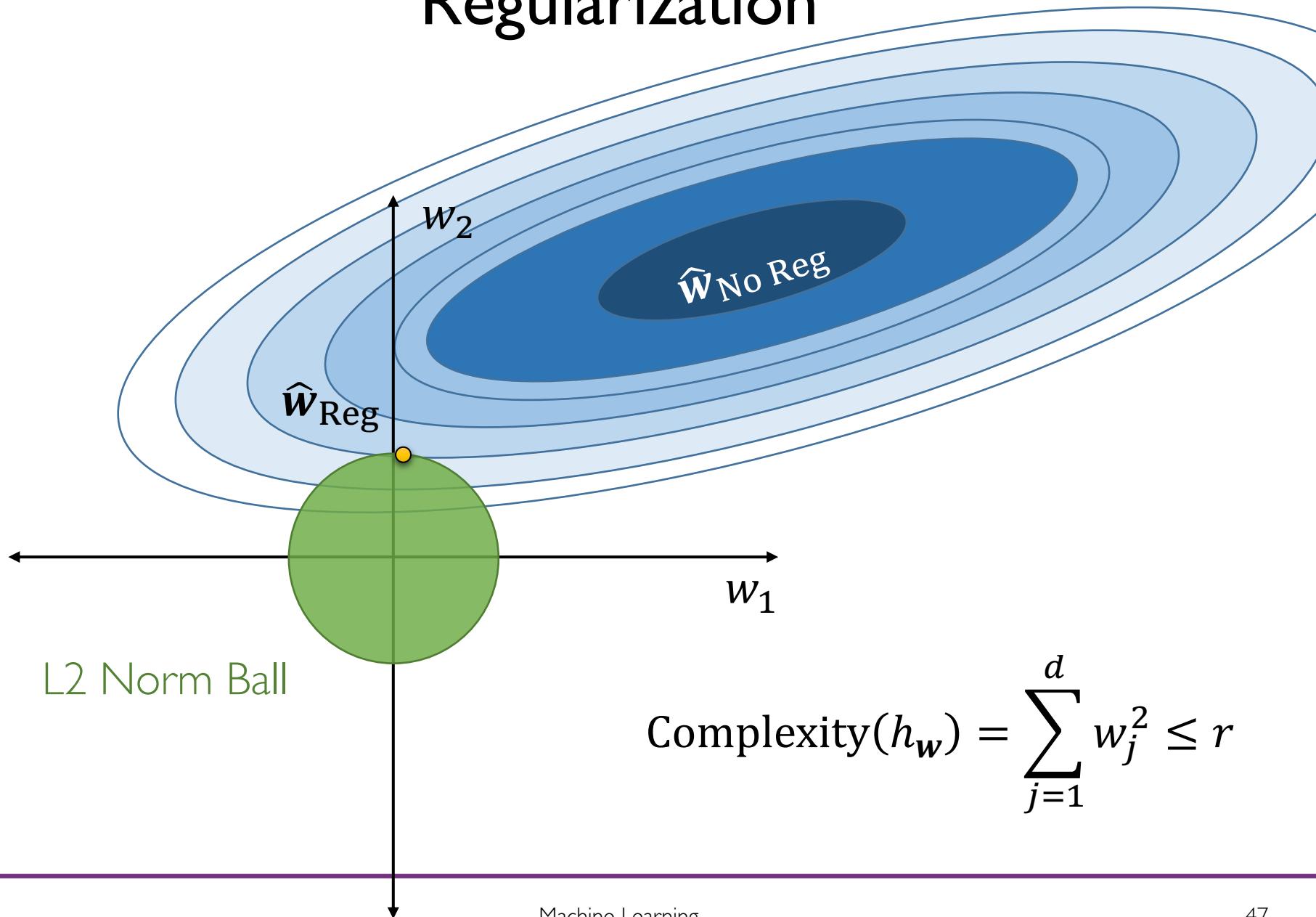
Regularization



Regularization

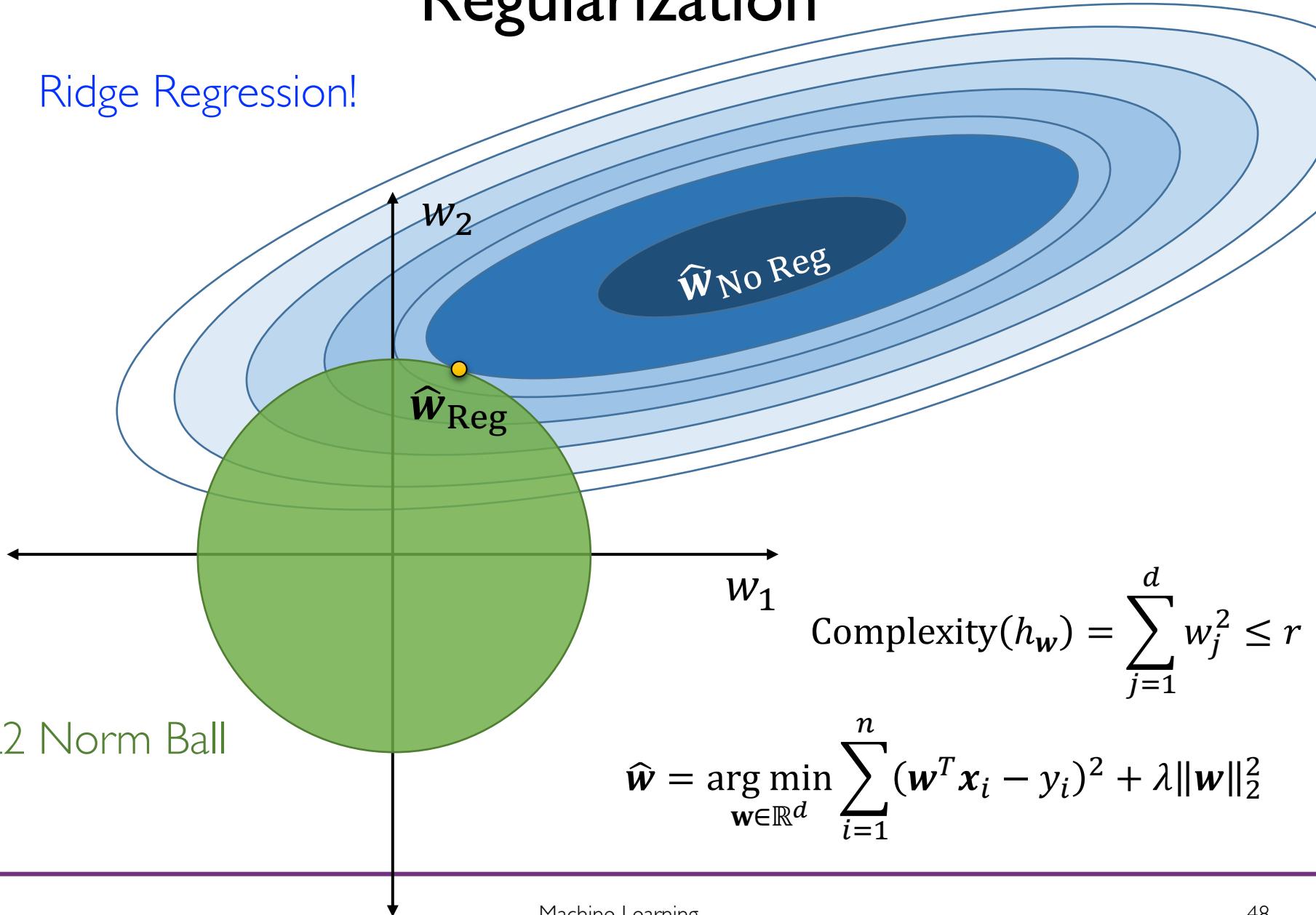


Regularization



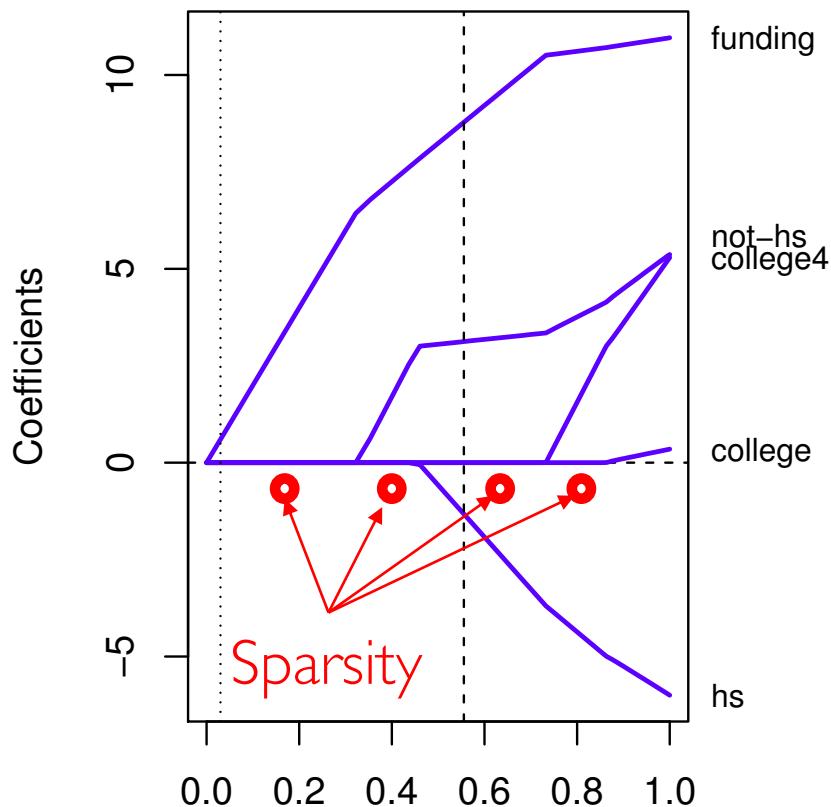
Regularization

Ridge Regression!



Regularization Path: L1 vs. L2

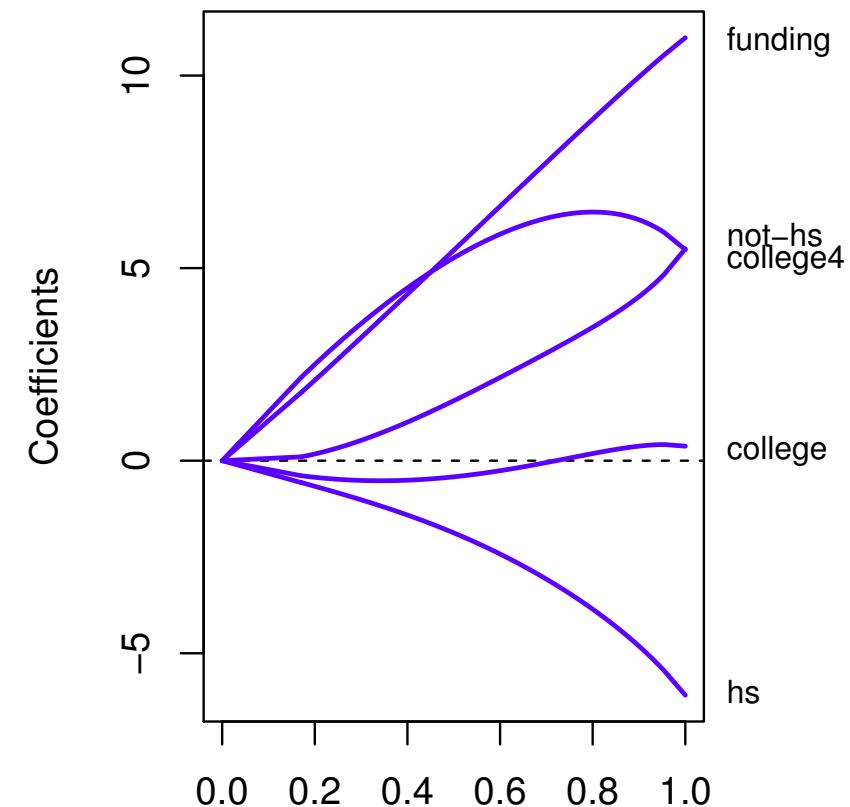
Lasso



$$\|\hat{\mathbf{w}}_r\|_1 / \|\hat{\mathbf{w}}\|_1$$

Feature selection

Ridge Regression



$$\|\hat{\mathbf{w}}_r\|_2 / \|\hat{\mathbf{w}}\|_2$$

Smooth optimization

$$\hat{\mathbf{w}}_r = \arg \min_{\|\mathbf{w}\|_1 \leq r} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$

Outline

- Linear Regression
 - Optimization
 - Statistical View: Maximum Likelihood Estimation
- Nonlinearization
- **Regularization**
 - L1 and L2 Norm
 - **Statistical View: Maximum a Posteriori Estimation**
- Linear Classification
 - Logistic Regression
 - Softmax Regression

Bayes Rule in Parametric Models

- Bayes rule:

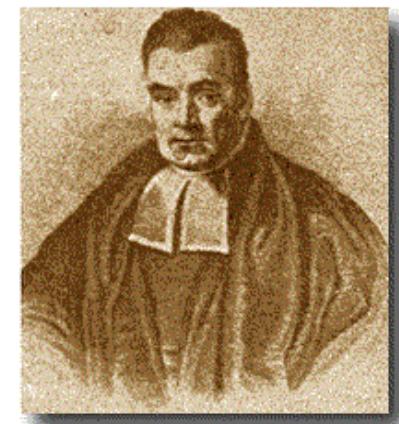
Observed Data Parameter

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

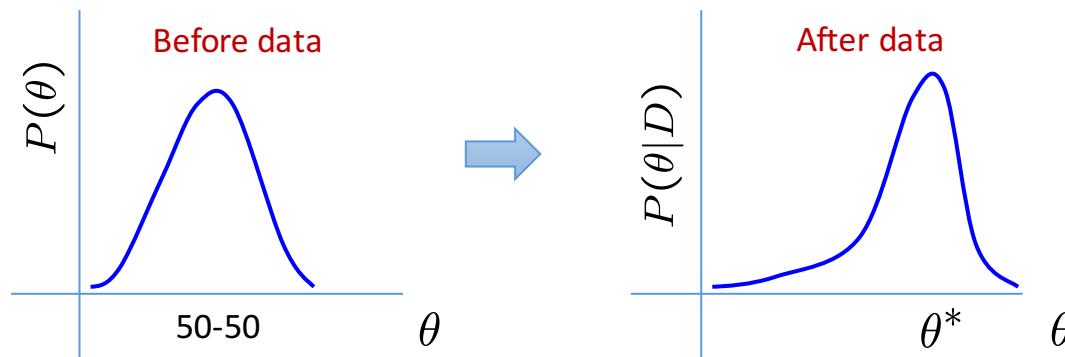
- Or equivalently:

$$p(\theta|D) \propto p(D|\theta)p(\theta)$$

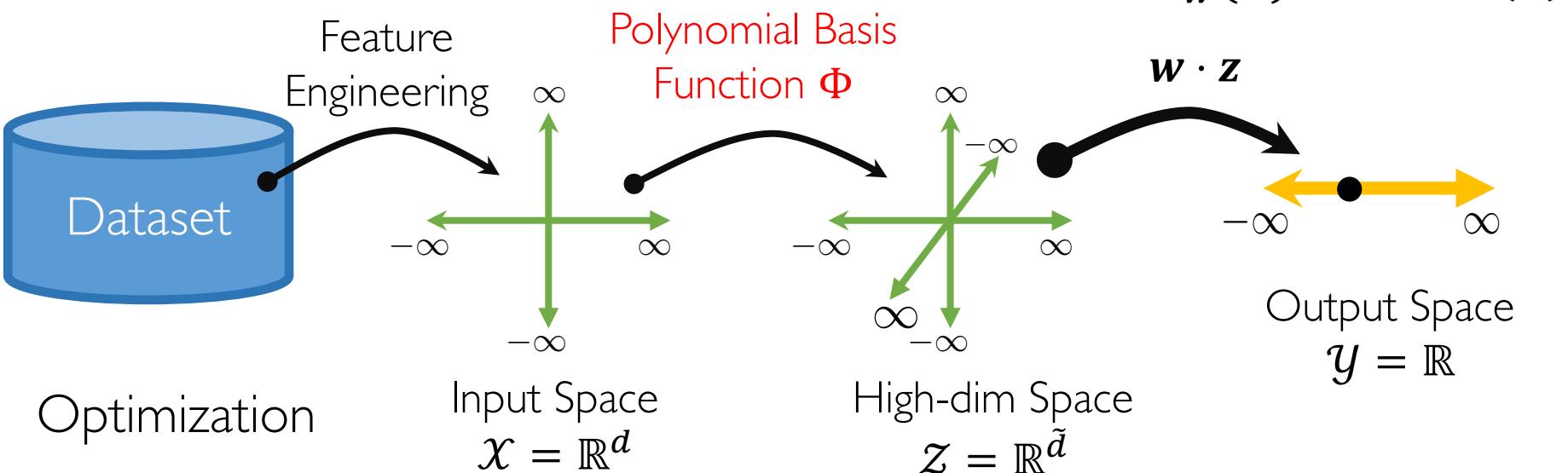
↓ ↓ ↓
Posterior Likelihood Prior



Thomas Bayes



Linear Regression



Optimization

Gradient Descent

$$\nabla_w \hat{\epsilon}(w) = 2\mathbf{Z}^T(\mathbf{Z}w - y) + 2\lambda w$$

Analytic Solution

$$w = (\mathbf{Z}^T \mathbf{Z} + \lambda I)^{-1} \mathbf{Z}^T y$$

Loss Function

$$\min_w \sum_{i=1}^n (h_w(x_i) - y_i)^2 + \lambda \Omega(w)$$

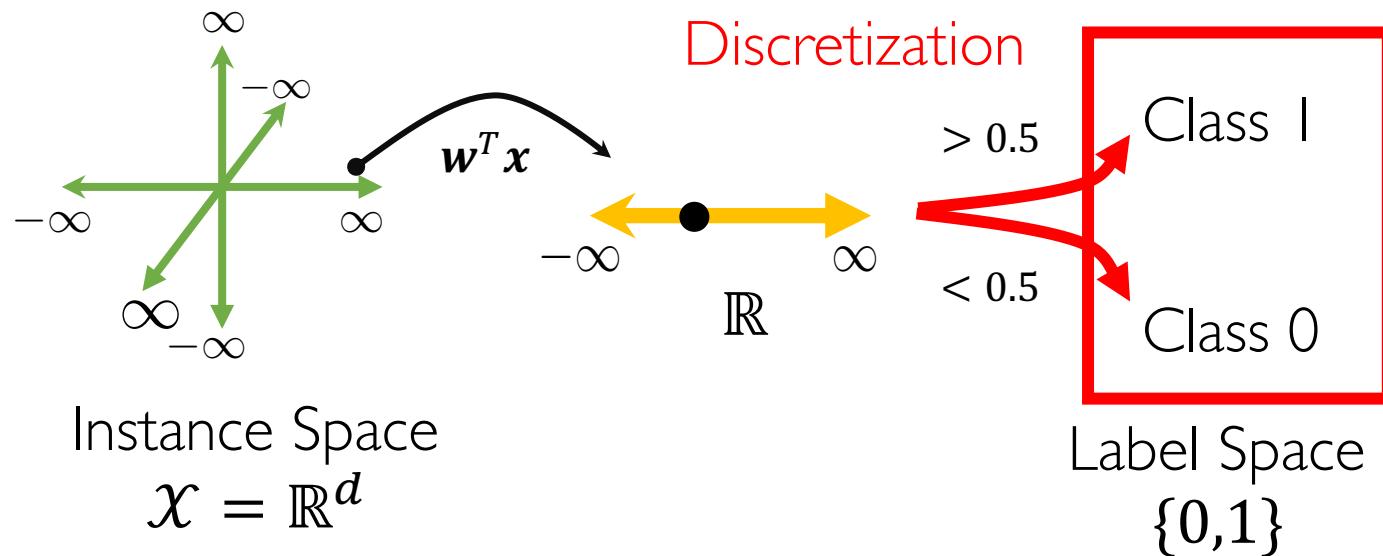
Well Done! A simple but complete model!

Outline

- Linear Regression
 - Optimization
 - Statistical View: Maximum Likelihood Estimation
- Nonlinearization
- Regularization
 - L1 and L2 Norm
 - Statistical View: Maximum a Posteriori Estimation
- **Linear Classification**
 - **Logistic Regression**
 - Softmax Regression

Regression to Classification

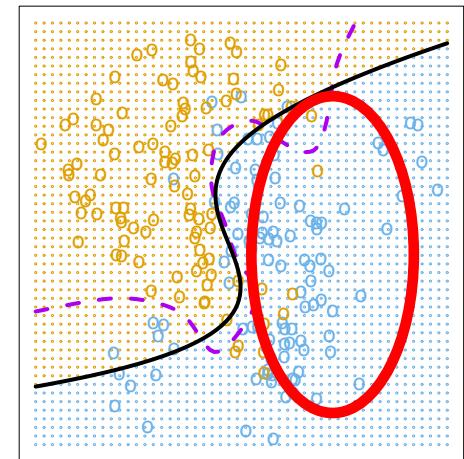
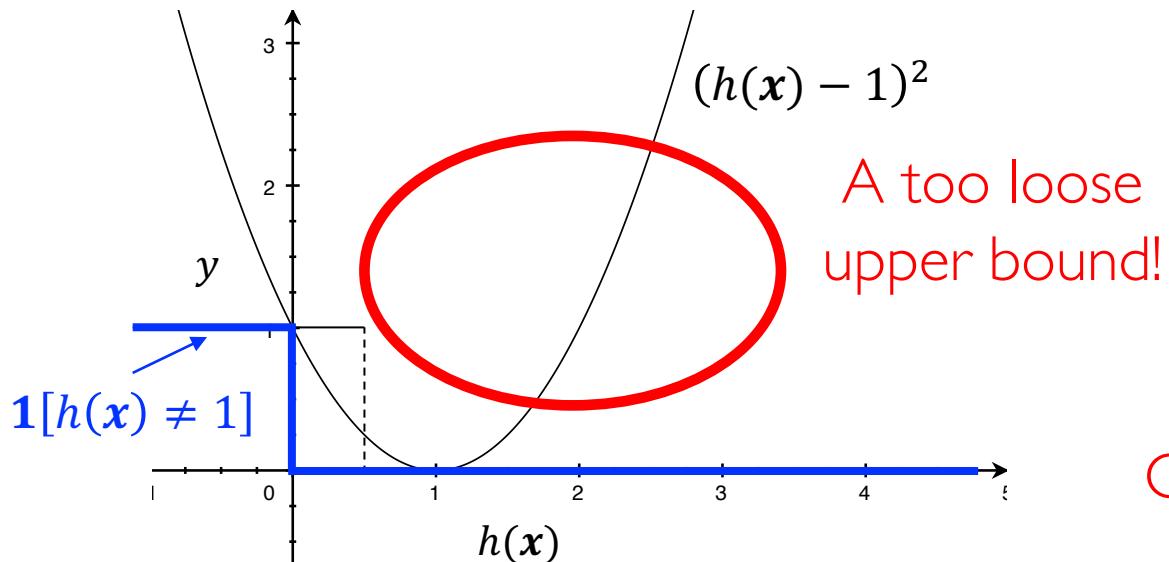
- How to design **loss function** over **linear hypothesis** for classification?
- Naïve idea: $\min_{\mathbf{w}} \sum_{i=1}^n (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 + \lambda \Omega(\mathbf{w})$ where $y_i \in \{0,1\}$
 - Discretize the continuous output $h_{\mathbf{w}}(\mathbf{x}_i)$ to be $\{0,1\}$
- During test part:



Regression to Classification

- It does not make sense for $h_{\mathbf{w}}$ to take values larger than 1 or smaller than 0 when we know that $y \in \{0,1\}$.
- We use accuracy (01-loss) when validating the classification model:

$$\ell(y, h(\mathbf{x})) = \mathbf{1}[y \neq h(\mathbf{x})]$$



- $h(\mathbf{x}) \in \mathcal{Y}$ is not a proper assumption this time.

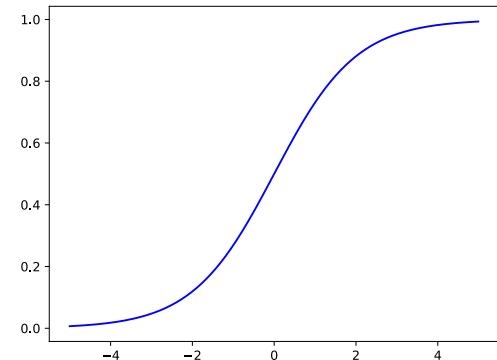
Logistic Regression: Link Function

- To map the output of $h(\mathbf{x})$ into $[0,1]$, we use **sigmoid function**:

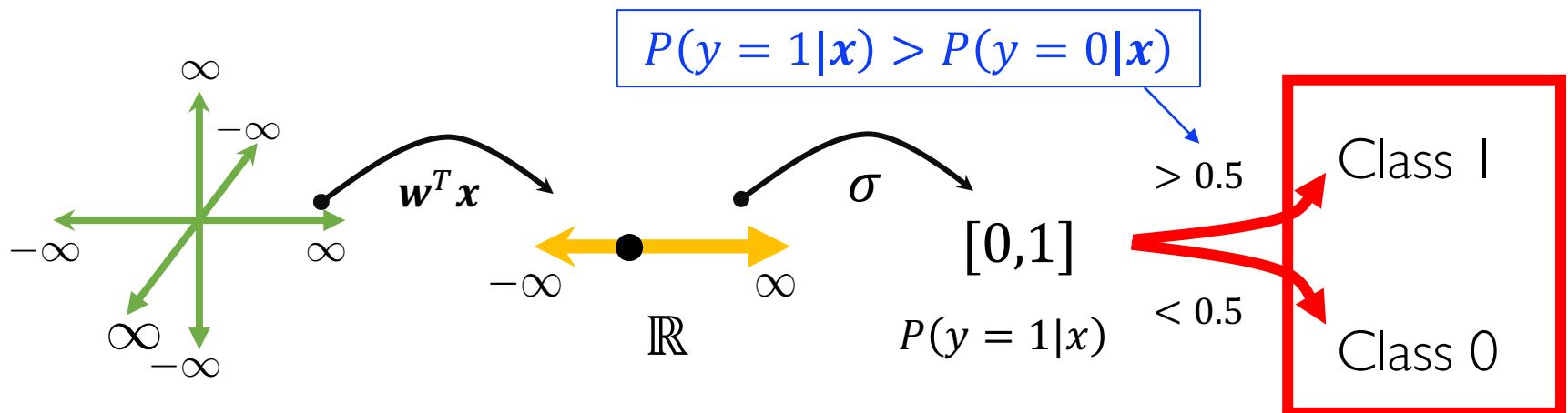
$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

- Sigmoid maps \mathbb{R} to **[0,1]**.

- $t \rightarrow +\infty, \sigma(t) \rightarrow 1$.



- Set $\sigma(h(\mathbf{x})) = p(y = 1|\mathbf{x})$ as the probability to label \mathbf{x} as $y = 1$.



Logistic Regression: Loss Function

- How about using **squared loss** now as $h(\mathbf{x}) \in [0,1]$ and $y \in \{0,1\}$?

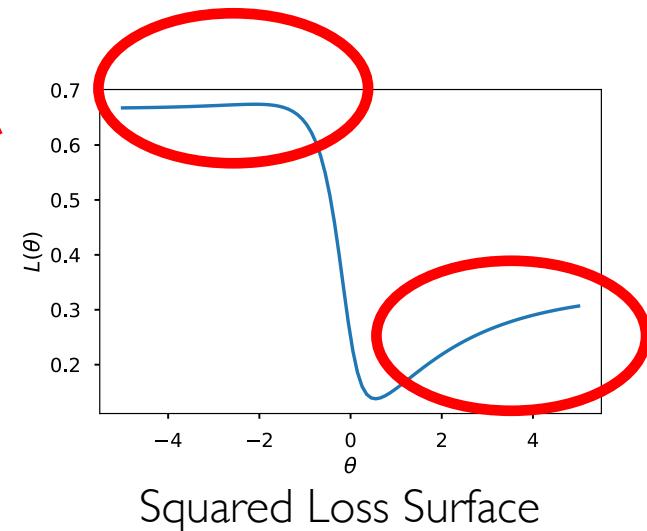
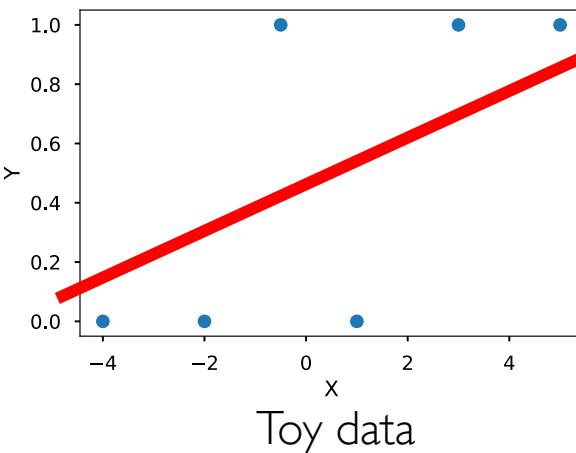
$$\ell(h(\mathbf{x}), y) = \sum_{i=1}^n (\sigma(\mathbf{w}^T \mathbf{x}_i) - y_i)^2 + \lambda \Omega(\mathbf{w})$$

- Tries to match **continuous** probability with **discrete** 0/1 labels.

- Non-convex

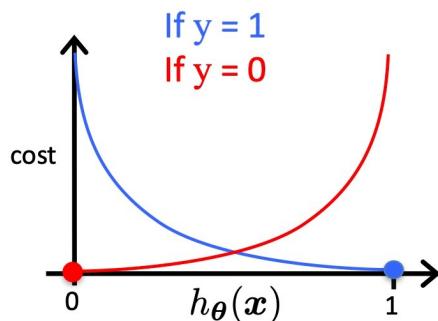
- Small loss when prediction is overly far in wrong side

- What is wrong?

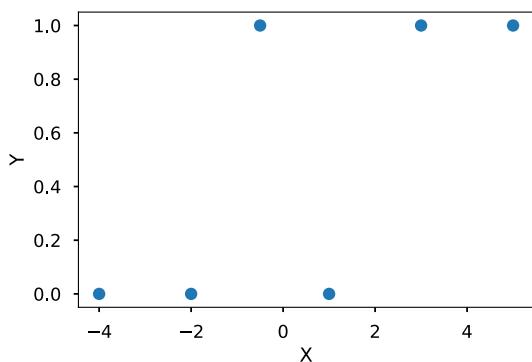


Logistic Regression: Cross-Entropy Loss

$$\ell(h(\mathbf{x}_i), y_i) = \begin{cases} -\log[\sigma(\mathbf{w}^T \mathbf{x}_i)] & y_i = 1 \\ -\log[1 - \sigma(\mathbf{w}^T \mathbf{x}_i)] & y_i = 0 \end{cases}$$



Toy data

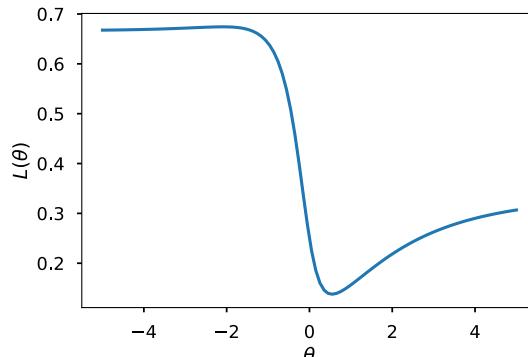


Check! If $y_i = 1$,

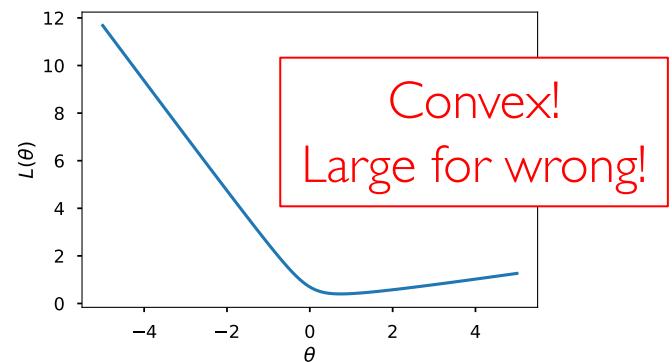
- $\sigma(\mathbf{w}^T \mathbf{x}_i) \rightarrow 0$, loss goes to ∞
- $\sigma(\mathbf{w}^T \mathbf{x}_i) \rightarrow 1$, loss goes to 0

$$P(y_i = 1 | \mathbf{x}_i)$$

Squared Loss Surface



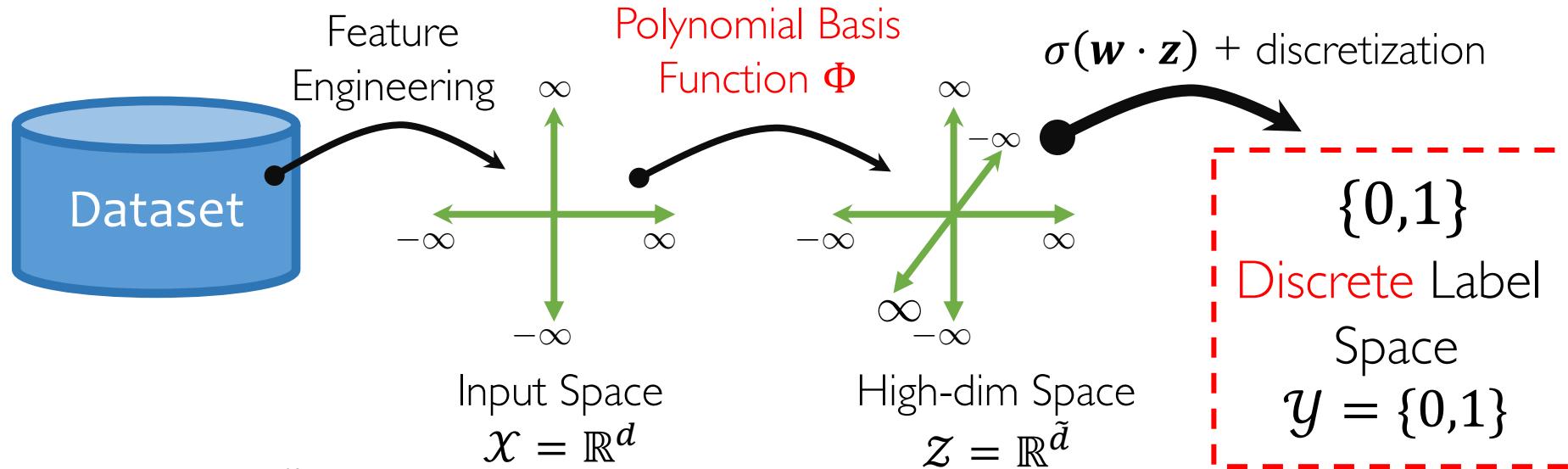
Cross-Entropy Surface



2-class
Classification:

Logistic Regression

Hypothesis Space
 $h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \Phi(\mathbf{x})$



$$\hat{\epsilon}(\mathbf{w}) = - \sum_{i=1}^n \{y_i \log \sigma(h_{\mathbf{w}}(\mathbf{x})) + (1 - y_i) \log[1 - \sigma(h_{\mathbf{w}}(\mathbf{x}))]\} + \lambda \Omega(\mathbf{w})$$

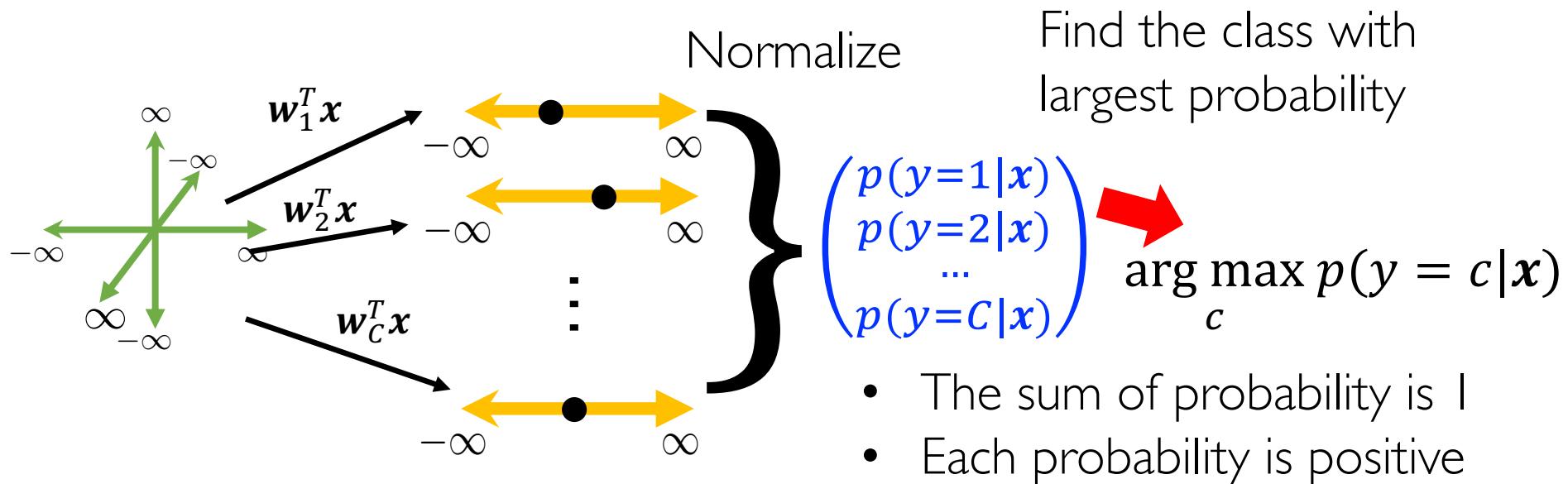
- How to solve this problem?
 - No analytic solution! Using GD, SGD, Newton method...
- How to deal with **multiclass** problems?

Outline

- Linear Regression
 - Optimization
 - Statistical View: Maximum Likelihood Estimation
- Nonlinearization
- Regularization
 - L1 and L2 Norm
 - Statistical View: Maximum a Posteriori Estimation
- **Linear Classification**
 - Logistic Regression
 - **Softmax Regression**

Multiclass Classification

- We keep on mapping the output of linear hypothesis into probability.
- There are multiple classes, so we use multiple linear predictors to predict the probability: $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C$.
- How to change the outputs of multiple predictors into probability?

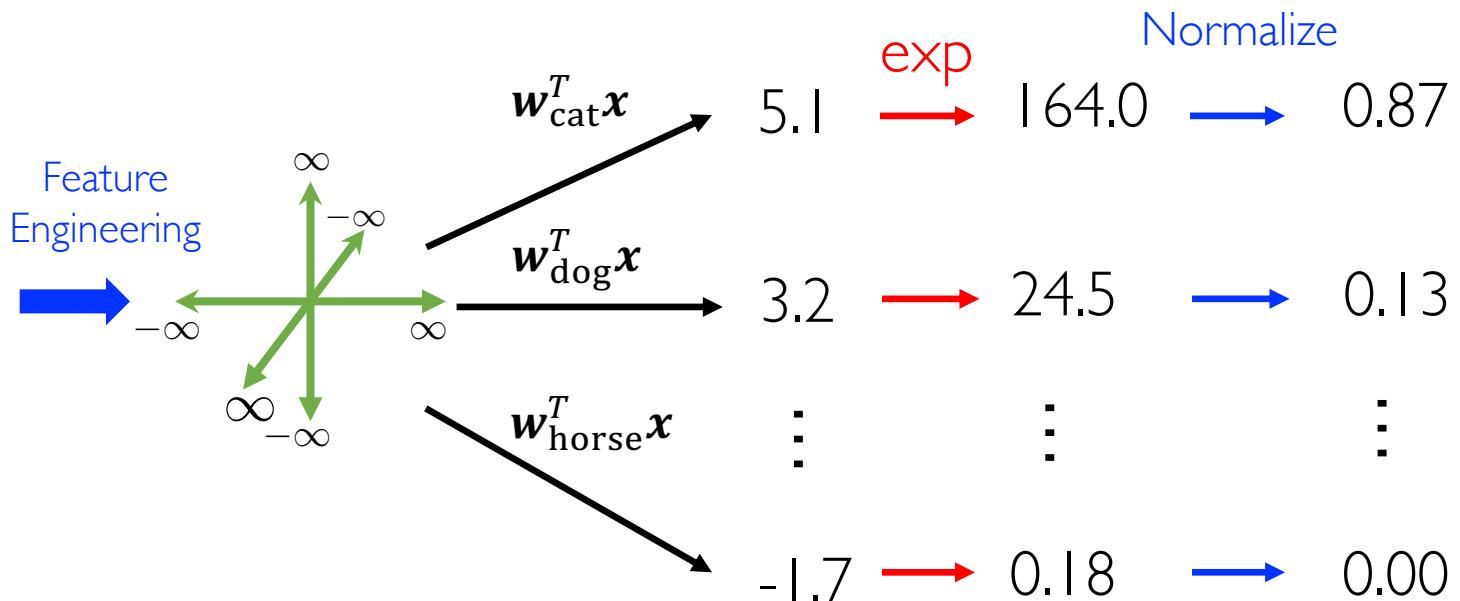


Softmax Regression: Softmax Function

- Softmax function normalizes multiple outputs in a probability vector:

$$p(y = c|x; \mathbf{W}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{r=1}^C \exp(\mathbf{w}_r^T \mathbf{x})}$$

Sum over
all classes



Softmax Regression: Cross-Entropy Loss

Loss for each data point (\mathbf{x}_i, y_i) :

Equivalent to maximum log-likelihood estimation

This loss is still convex. 

Minimize this loss: GD, SGD...

$$\ell(h(\mathbf{x}_i), y_i) = \begin{cases} -\log \frac{\exp(\mathbf{w}_1^T \mathbf{x}_i)}{\sum_{r=1}^C \exp(\mathbf{w}_r^T \mathbf{x}_i)}, & y_i = 1 \\ -\log \frac{\exp(\mathbf{w}_2^T \mathbf{x}_i)}{\sum_{r=1}^C \exp(\mathbf{w}_r^T \mathbf{x}_i)}, & y_i = 2 \\ \vdots \\ \vdots \\ -\log \frac{\exp(\mathbf{w}_C^T \mathbf{x}_i)}{\sum_{r=1}^C \exp(\mathbf{w}_r^T \mathbf{x}_i)}, & y_i = C \end{cases}$$

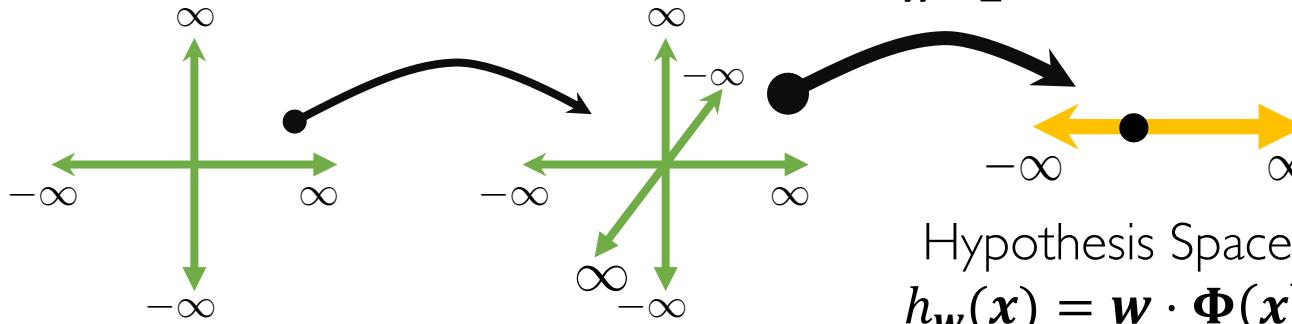
These two hypotheses have same outputs!

$$\begin{pmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \vdots \\ \mathbf{w}_C \end{pmatrix} \quad \begin{pmatrix} \mathbf{w}_1 + V \\ \mathbf{w}_2 + V \\ \vdots \\ \mathbf{w}_C + V \end{pmatrix}$$

Add regularization!
Or the parameter will go to infinity.

Linear Model: Summary

I. Basis Function Φ



Hypothesis Space
 $h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \Phi(\mathbf{x})$

Objective function
may not be linear to
the parameters.

4. Optimization

- High-dimension
 - GD
 - SGD
- High-dimension big data
- Non-differentiable
 - Subgradient

$$\min_{\mathbf{w}} \sum_{i=1}^n \ell(h_{\mathbf{w}}(\mathbf{x}_i), y_i) + \lambda \Omega(\mathbf{w})$$

2. Loss Function

- Regression
 - Squared Loss
- Classification
 - Cross-Entropy

3. Regularization

- General
 - L2 Regularizer
- Sparse Solution
 - L1 Regularizer