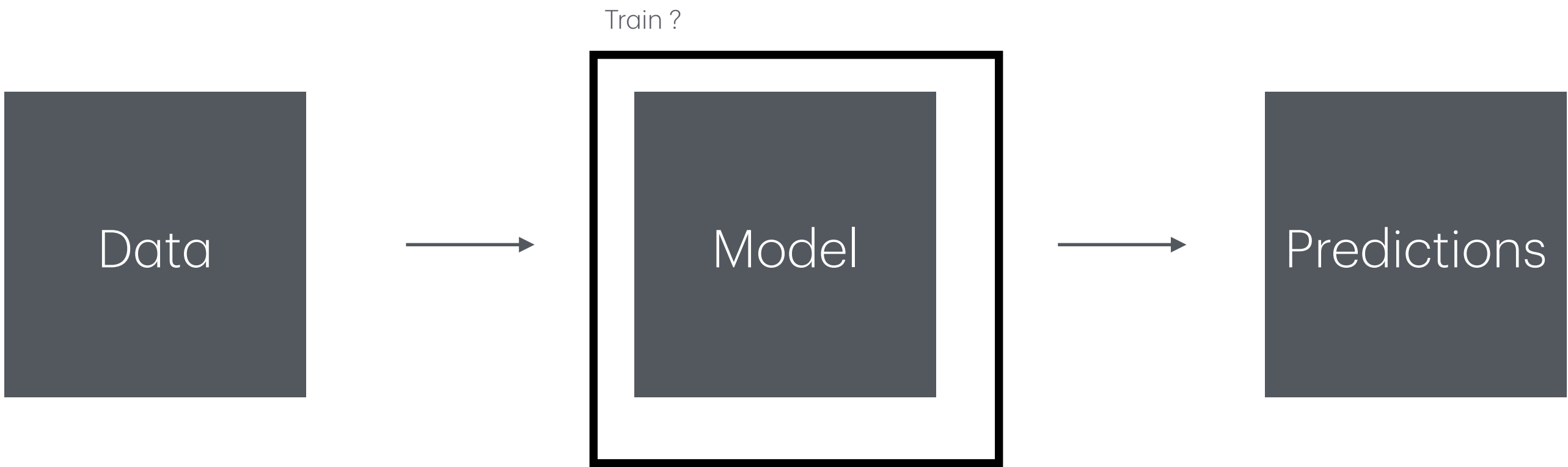# Training Models

# Polynomial Models

Prone to overfitting. Use this as an guide to handle real world problems

- How to detect overfitting

- Use regularization to reduce the risk of overfitting

# Linear Regression

Linear based model: $\theta_0 + \theta_1 x_1 + \theta_2 x_2 \dots$

Train ?



Data → Model → Predictions

Train model by finding models parameters that lower cost function.

Cost function is the sum of MSE of each instance prediction

Goal: Find $\theta$ vector which minimizes prediction MSE

| X (Instance) | h(X) (Prediction) | Y (Target) | Error ( h(X) - Y ) | Error^2 ( h(X) - Y )^2 |
|---|---|---|---|---|
|  |  |  |  |  |
|  |  |  |  |  |
|  |  |  |  |  |
|  |  |  |  |  |

MSE Error

$$\frac{1}{m}\Sigma Error_{m_i}(\dots)^2$$

# MSE Error

Normal Equation
(Closed Form Solution)

Gradient Descent

# Normal Equation: MSE Error

Option 1 uses SVD(Singular Value Decomposition) algorithm

Normal Equation
(Closed Form Solution)

Option 1

$$\hat{\theta} = X^+ \cdot y$$

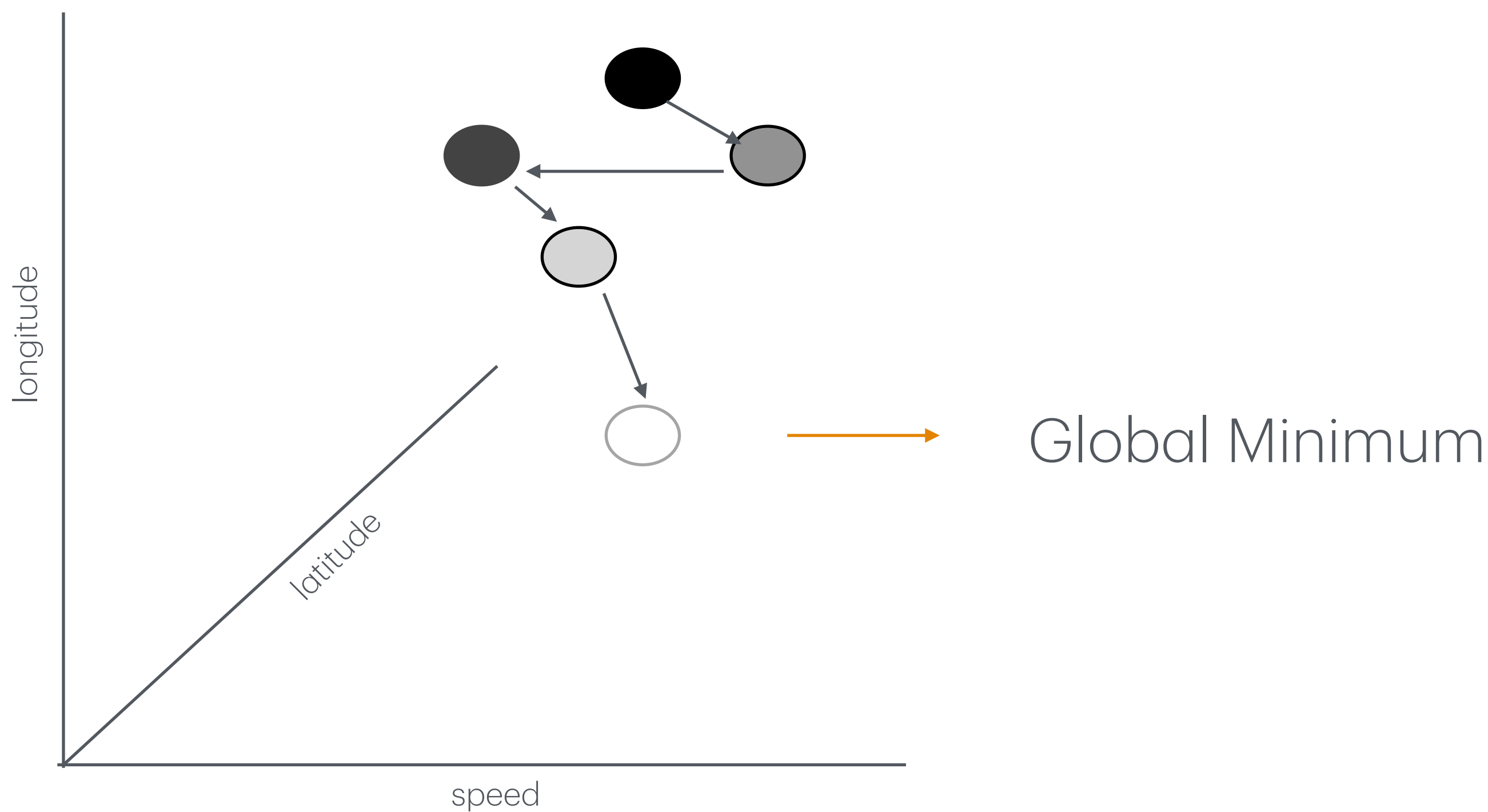$$X^+ = \text{Pseudoinverse}$$

Option 2

$$\hat{\theta} = (X^T \cdot X)^{-1} \cdot X^T \cdot y$$
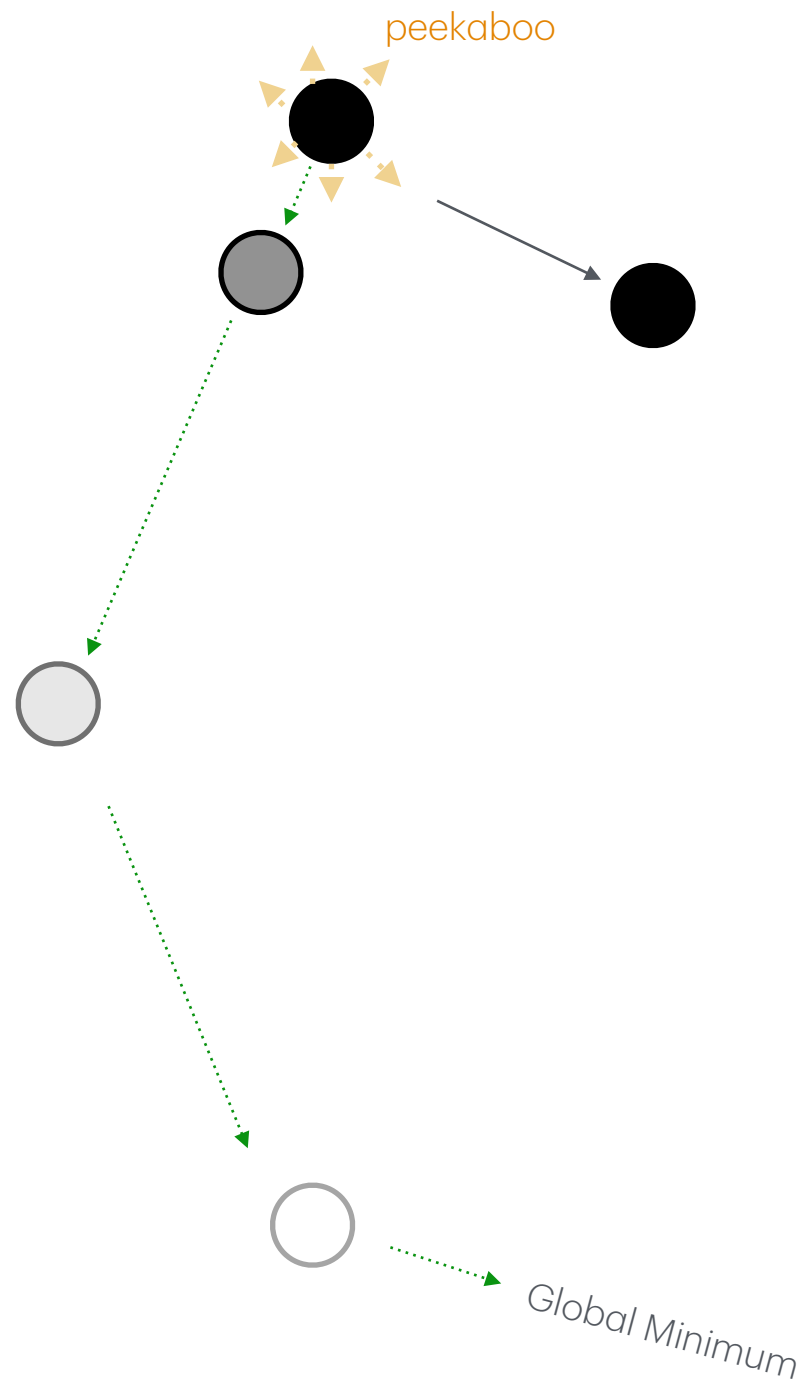
Both options are slow for very large datasets

# Gradient Descent: MSE Error

$$X_{speed} + X_{longitude} + X_{latitude}$$

$$\theta_{speed} + \theta_{longitude} + \theta_{latitude} \longrightarrow$$

Find solution which minimizes error



Global Minimum

# Gradient Descent: MSE Error

peekaboo

Global Minimum

Batch Gradient Descent

$$X = X_{i_{speed}} + X_{i_{longitude}} X_{i_{latitude}}$$

$$\hat{\theta} = \theta_{speed} + \theta_{longitude} + \theta_{latitude}$$

| i | speed | longitude | latitude | Y |
|---|-------|-----------|----------|---|
| 0 |       |           |          |   |
| 1 |       |           |          |   |

$\vdots$

$$MSE(X_i, \hat{\theta})$$

$$\nabla = \frac{\partial}{\partial \theta_{speed}} + \frac{\partial}{\partial \theta_{longitude}} + \frac{\partial}{\partial \theta_{latitude}}$$

$$\nabla MSE(X_i, \hat{\theta})$$

Gradient of MSE (Cost Function)

$$\frac{\partial}{\partial \theta_{speed}} \cdot MSE(X_i, \hat{\theta}) \longrightarrow$$

$$\frac{\partial}{\partial \theta_{longitude}} \cdot MSE(X_i, \hat{\theta}) \longrightarrow$$

$$\frac{\partial}{\partial \theta_{latitude}} \cdot MSE(X_i, \hat{\theta}) \longrightarrow$$
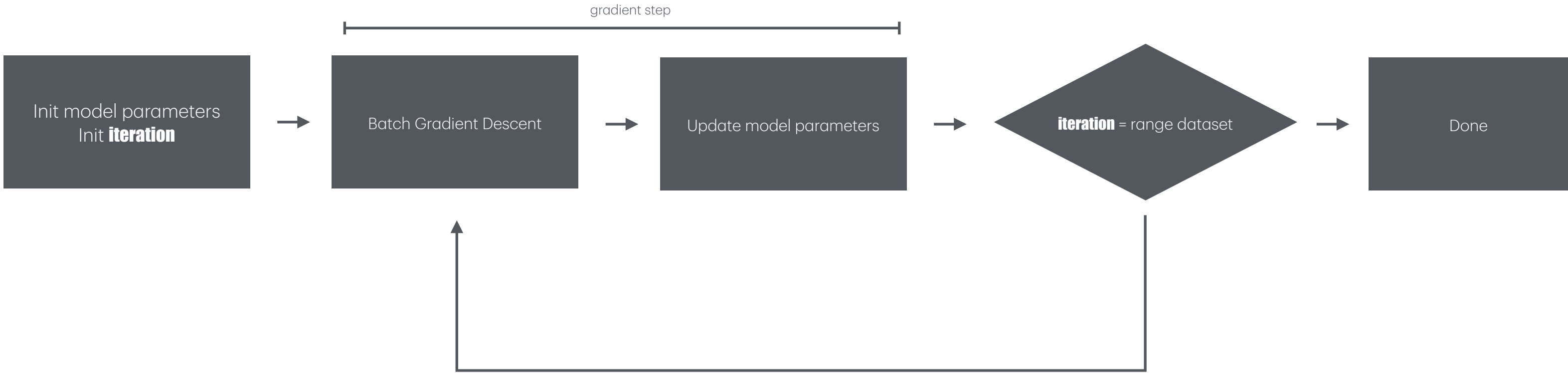
$$= \frac{2}{m} X^T (X\theta - y)$$

Requires **entire training dataset (Batch Gradient Descent)** per training step,

Vector with length equal to number of features

Scales with with feature (i.e. preferred algorithm for large number of features)
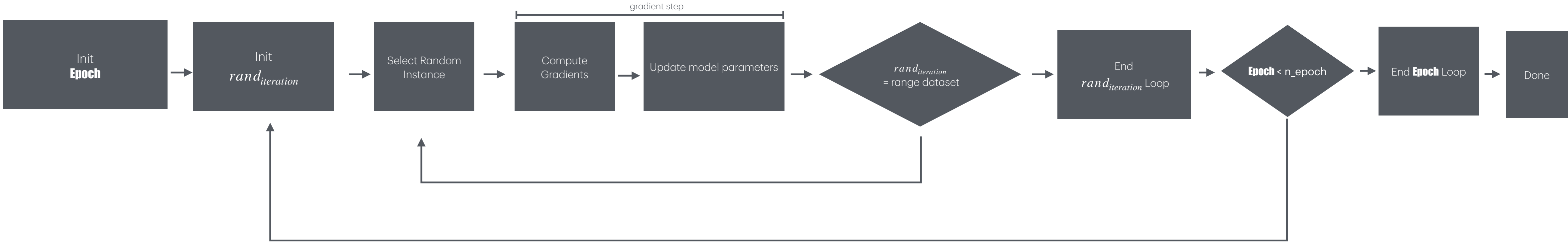
# Batch Gradient Descent

gradient step

| Init model parameters Init **iteration** | → | Batch Gradient Descent | → | Update model parameters | → | **iteration** = range dataset | → | Done |

# Stochastic Gradient Descent

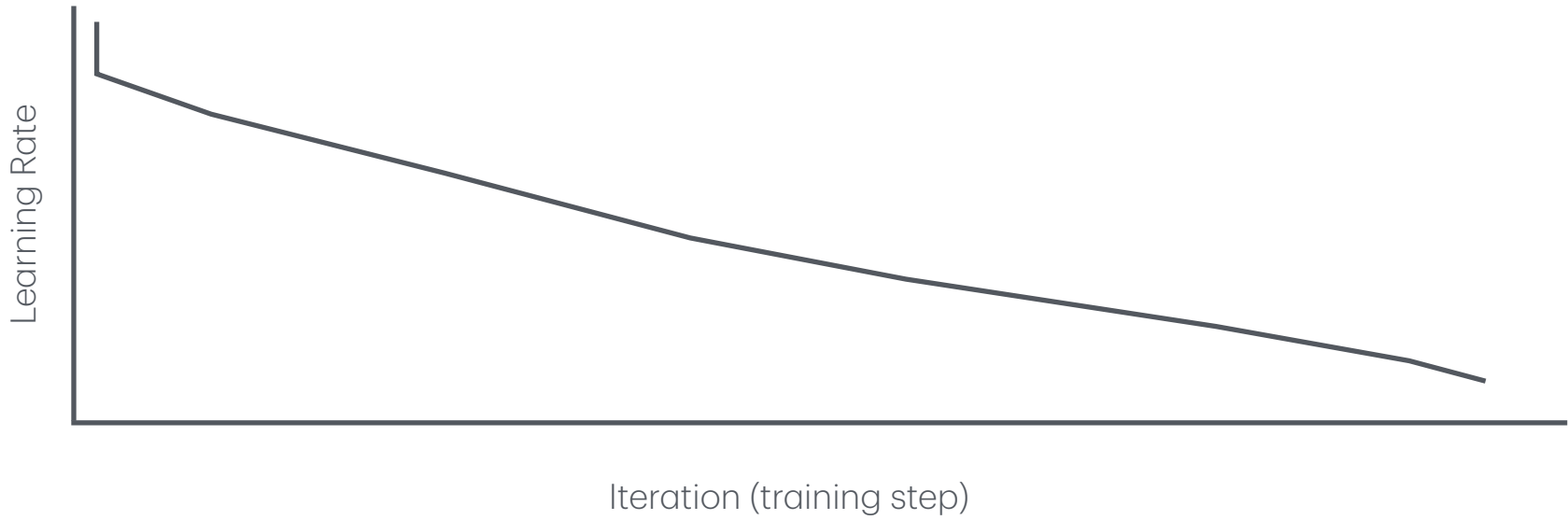| i | random i | speed | longitude | latitude | y |
|---|---|---|---|---|---|
| 0 | 2 | | | | |
| 1 | 6 | | | | |
| 2 | | | | | |
| 3 | 4 | | | | |
| 4 | 7 | | | | |
| 5 | 1 | | | | |
| 6 | 3 | | | | |
| 7 | | | | | |
| 8 | | | | | |

gradient step

**Init Epoch** → **Init $rand_{iteration}$** → **Select Random Instance** → **Compute Gradients** → **Update model parameters** → **$rand_{iteration}$ = range dataset** → **End $rand_{iteration}$ Loop** → **Epoch < n_epoch** → **End Epoch Loop** → **Done**
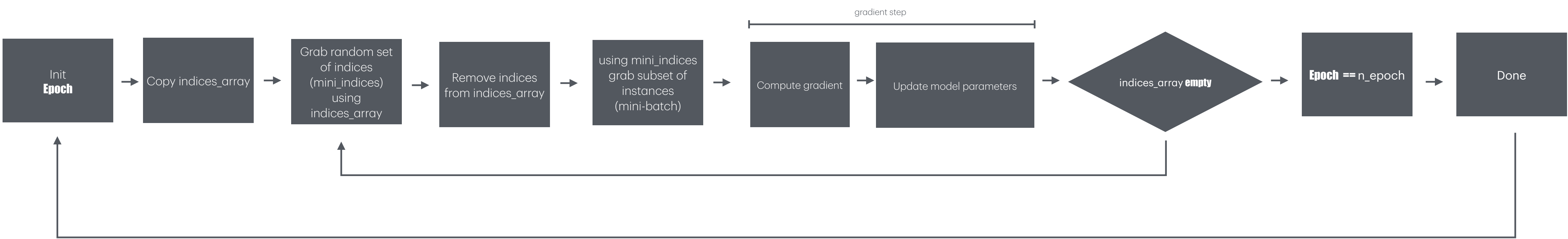
Each training step is much faster but the search for the minimum is noisy (bouncy, random)

Algorithm alone does not find true minimum. Gradually reduce learning rate during training (simulated annealing)

Learning Rate

Iteration (training step)

Mini-Batch Gradient Descent

gradient step

| Init **Epoch** | → | Copy indices_array | → | Grab random set of indices (mini_indices) using indices_array | → | Remove indices from indices_array | → | using mini_indices grab subset of instances (mini-batch) | → | Compute gradient | Update model parameters | → | indices_array **empty** | → | **Epoch** == n_epoch | → | Done |

Algorithm provides a performance boost from hardware optimization of matrix operations (perfect for GPU matrix processing)
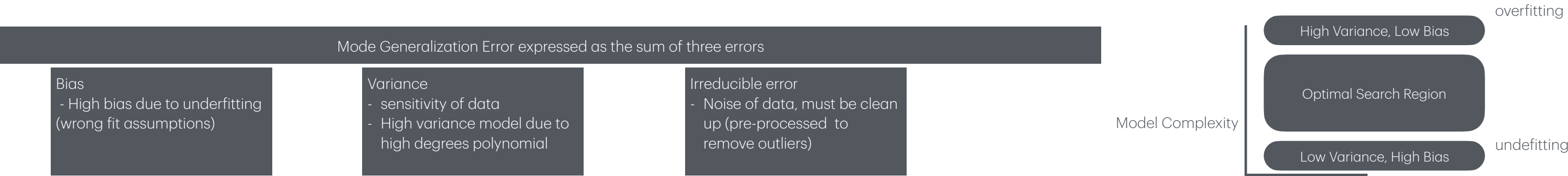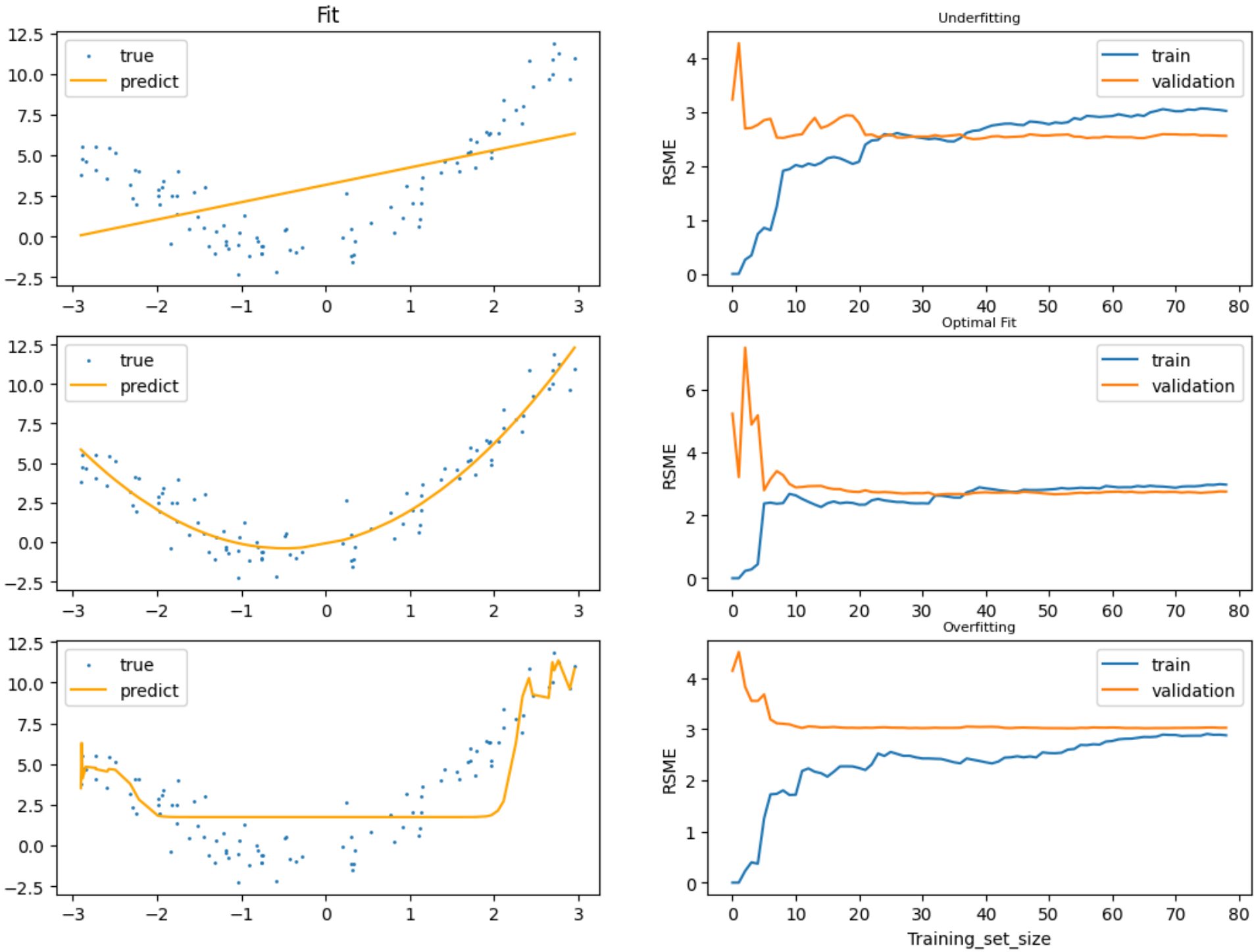
Prone to finding local minima if the learning schedule is bad

# Old Way to measure model

- Training performs well, Test set does not
  ‣ Overfitting
- Training and Test performs bad
  ‣ Underfitting

# New Way to measure model

- Use Learning Curves
  - Using Training Set
  - Using Validation Set

Data ( no idea what/who generated )

Model

How complex should it be?

Measures overfitting and underfitting?



Mode Generalization Error expressed as the sum of three errors

| Bias | Variance | Irreducible error |
|---|---|---|
| - High bias due to underfitting (wrong fit assumptions) | - sensitivity of data - High variance model due to high degrees polynomial | - Noise of data, must be clean up (pre-processed to remove outliers) |

overfitting

High Variance, Low Bias

Optimal Search Region

Model Complexity

Low Variance, High Bias

undefitting

Linear based model: $\theta_0 + \theta_1 x_1 + \theta_2 x_2 \ldots$

Constrain the weight of the model

| X1 (Instance) | X2 (Instance) | h(X) (Prediction) | Y (Target) | Error ( h(X) - Y ) | Error^2 ( h(X) - Y )^2 |
|---|---|---|---|---|---|
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

$MSE(\theta_i) +$ Ridge-Regularizer

$$Cost = MSE(\theta_i) + \alpha \frac{1}{2}(\theta_1^2 + \theta_2^2)$$

Use **ONLY** during training (updated model parameters)
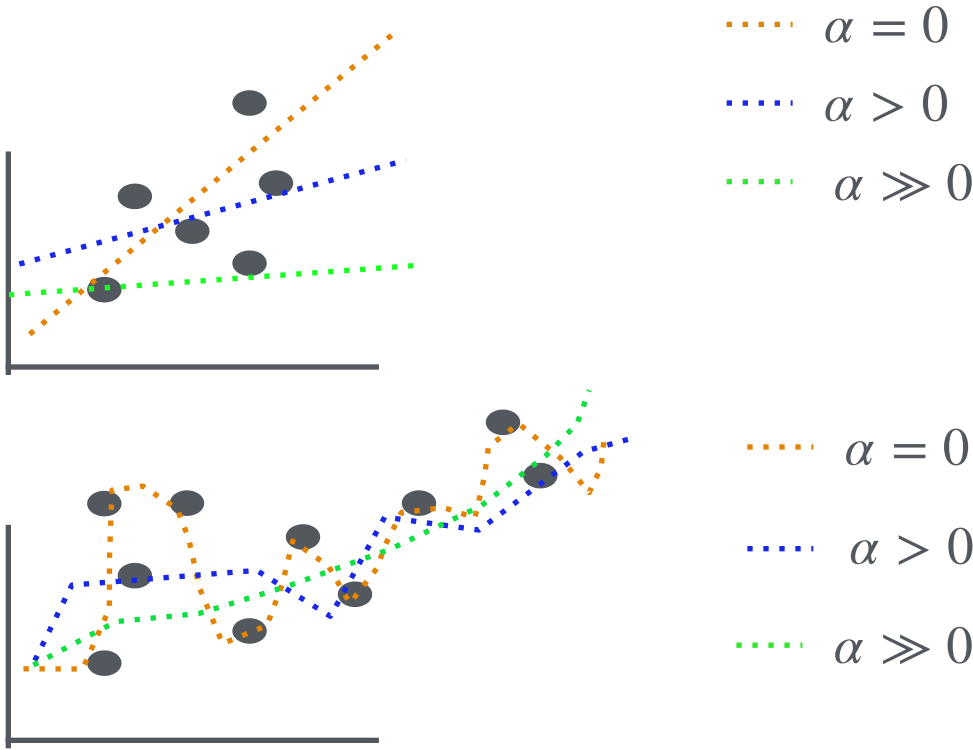
$$Cost = MSE(\theta_i)$$

Use during evaluation

$$Cost = MSE(\theta_i) + \alpha\frac{1}{2}(\theta_1^2 + \theta_2^2)$$

Training step works to reduce MSE

Training step works to reduce model weights

Reduces slope of 1 degree fit

Reduces sensitivity of 2 degree fit

$\alpha = 0$
$\alpha > 0$
$\alpha \gg 0$

$\alpha = 0$
$\alpha > 0$
$\alpha \gg 0$

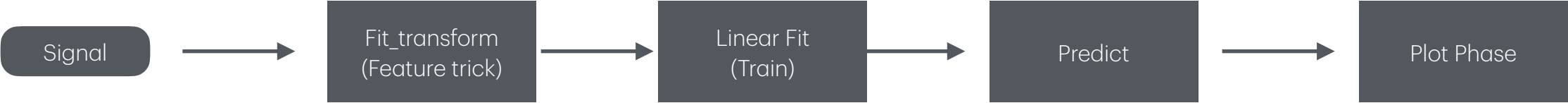Regularization reduces generalization error's variance(increases bias)

$$Cost = MSE(\theta_i) + \alpha \frac{1}{2}(\theta_1^2 + \theta_2^2)$$
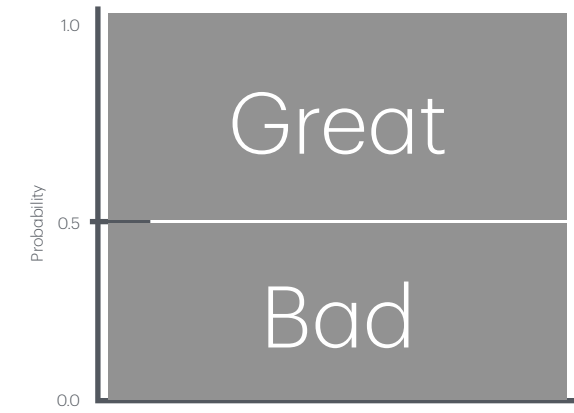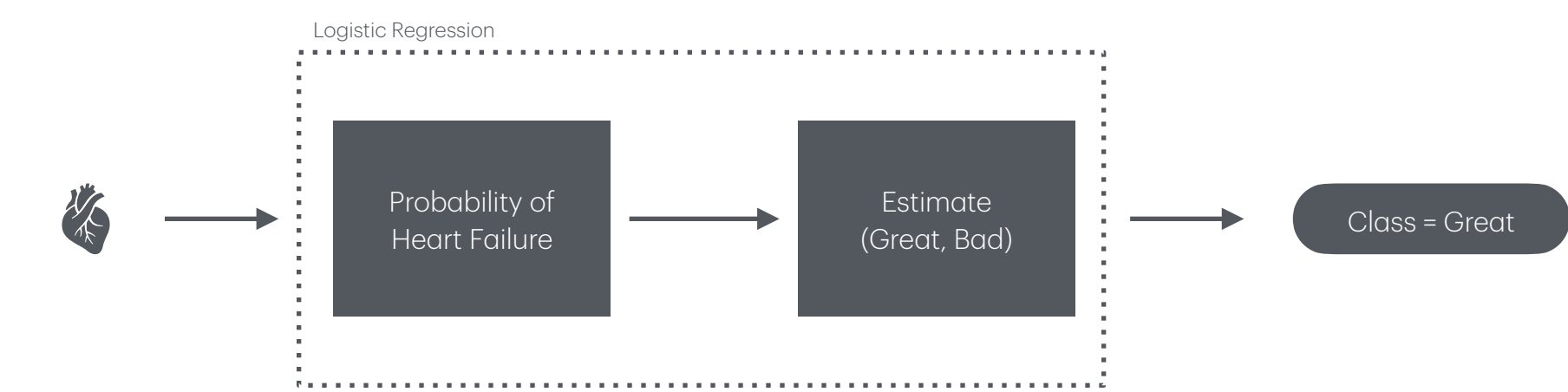
Training step
works to
reduce MSE

Ridge
Regularization

Lasso Regularization
another regularizer which
eliminates the weights(i.e. $\theta_i$)
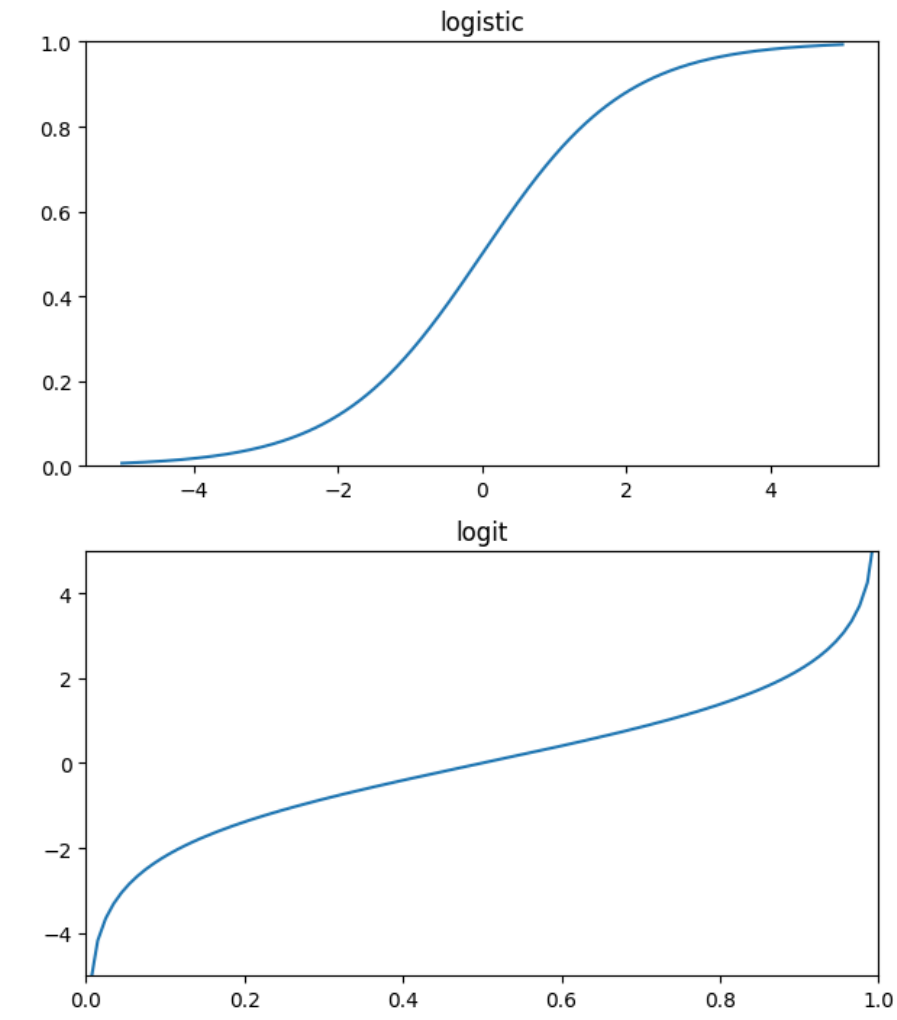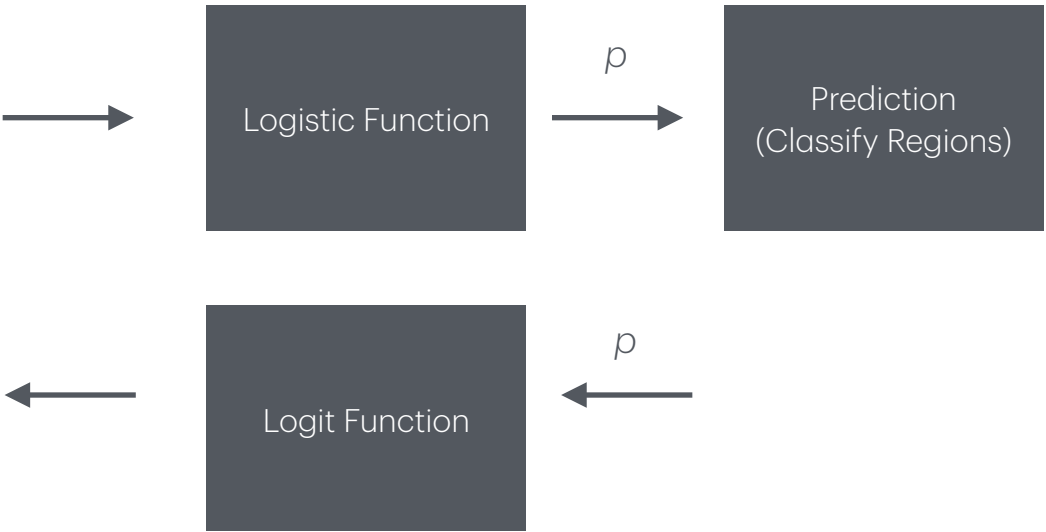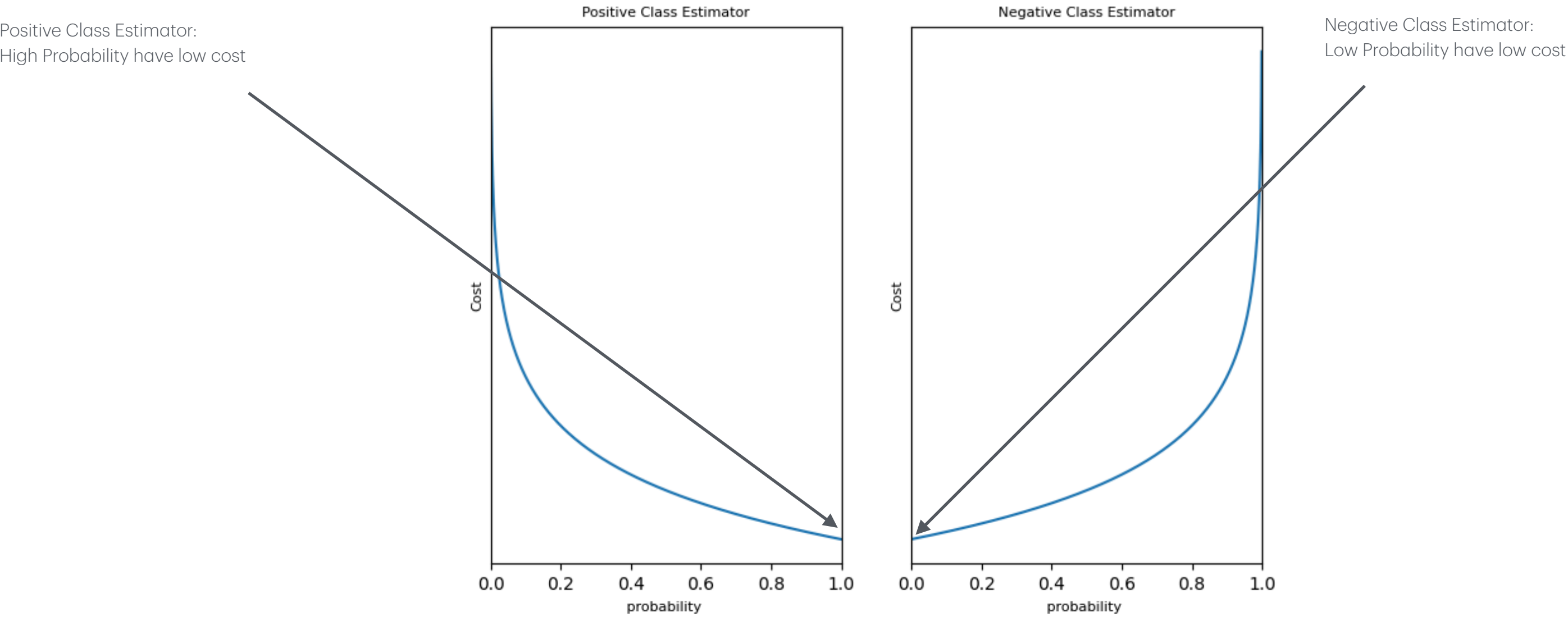of least important features

# Pipeline Poly

Signal → Fit_transform (Feature trick) → Linear Fit (Train) → Predict → Plot Phase

# Logistic Regression

Logistic Regression

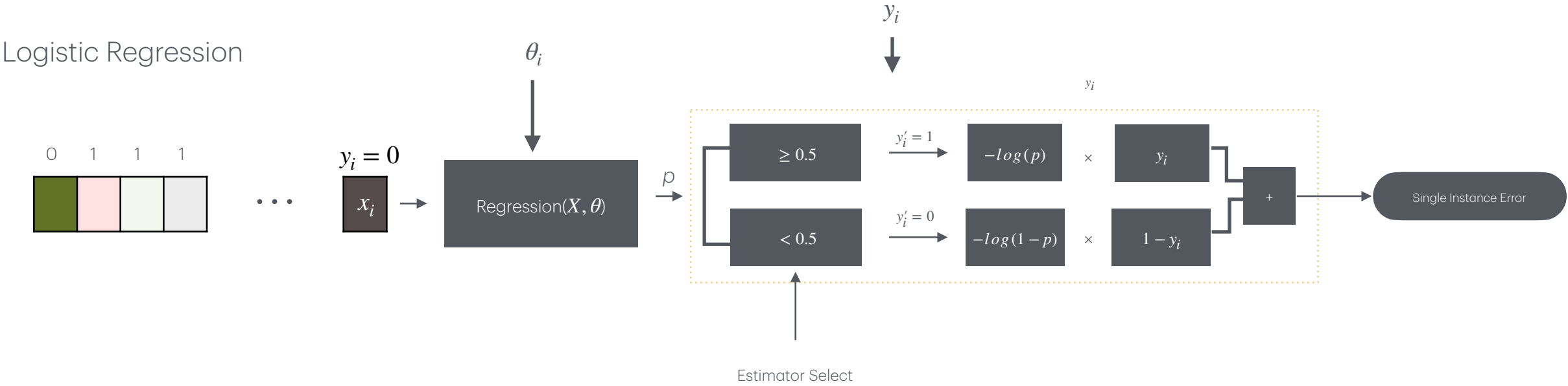Probability of Heart Failure → Estimate (Great, Bad) → Class = Great



Great

Bad

Probability

Weighted Sum: $\theta_0 + \theta_1 x_1 + \theta_2 x_2 \ldots$ →

Logistic Function → $p$ → Prediction (Classify Regions)

Logit Function ← $p$



logistic

logit

# Logistic Regression: How to Train?

| Positive Class | → Parameters($\theta_i$) estimates high probability for positive instance |

| Negative Class | → Parameters estimates low probability for negative instance |

### Linear Regression

$\theta_i$     $y_i$

| 0 | 1 | 1 | 1 | $\cdots$ | $y_i = 0$ $x_i$ → | Regression($X, \theta$) → | MSE → | Single Instance Error |

### Logistic Regression

$\theta_i$     $y_i$

| 0 | 1 | 1 | 1 | $\cdots$ | $y_i = 0$ $x_i$ → | Regression($X, \theta$) |

$p$

$\geq 0.5$   $y_i' = 1$   $-log(p)$   $\times$   $y_i$

$< 0.5$   $y_i' = 0$   $-log(1-p)$   $\times$   $1 - y_i$

$+$ → Single Instance Error

Estimator Select

$\theta_i$

Positive Class Estimator:
High Probability have low cost

Negative Class Estimator:
Low Probability have low cost



Positive Class Estimator

Negative Class Estimator
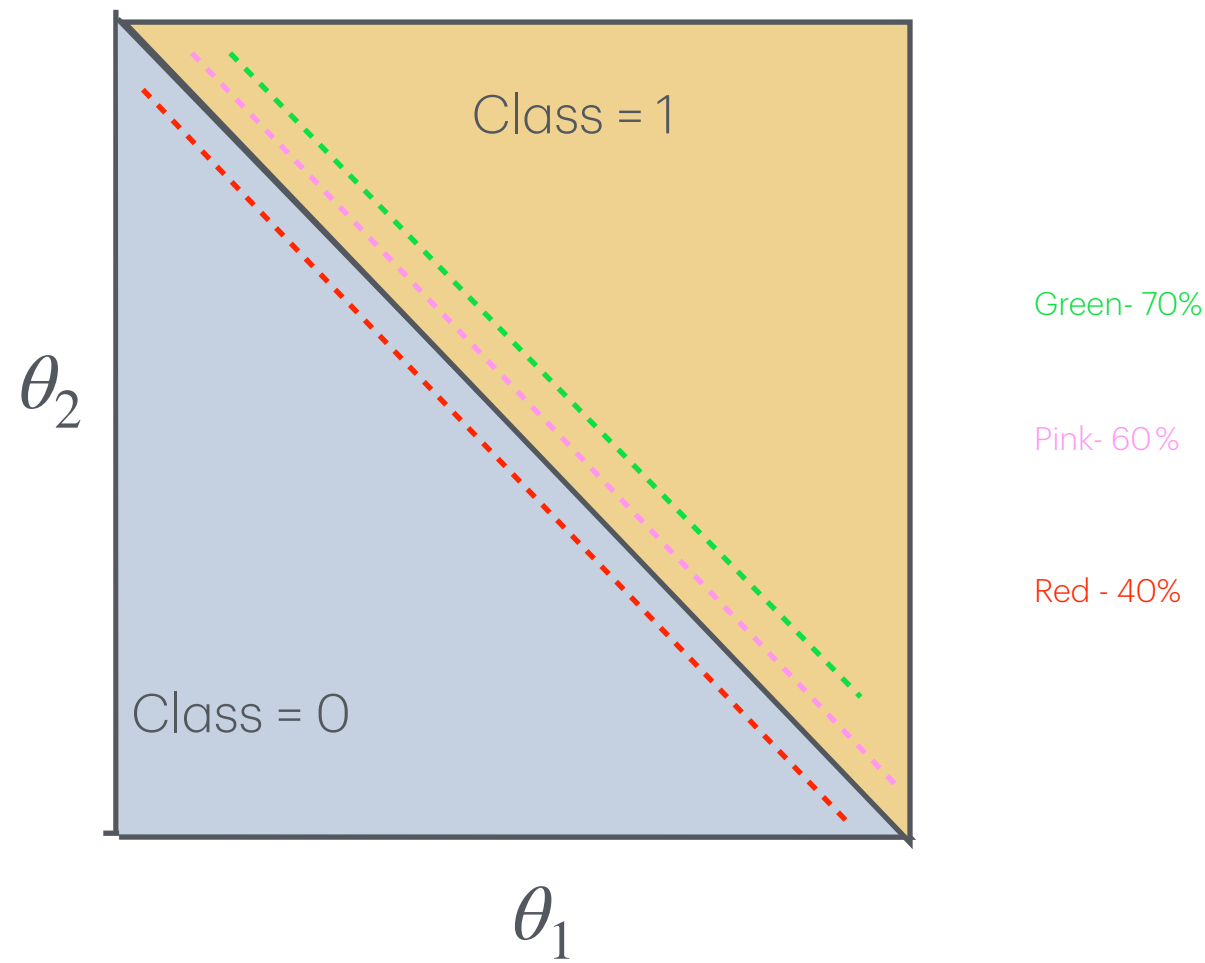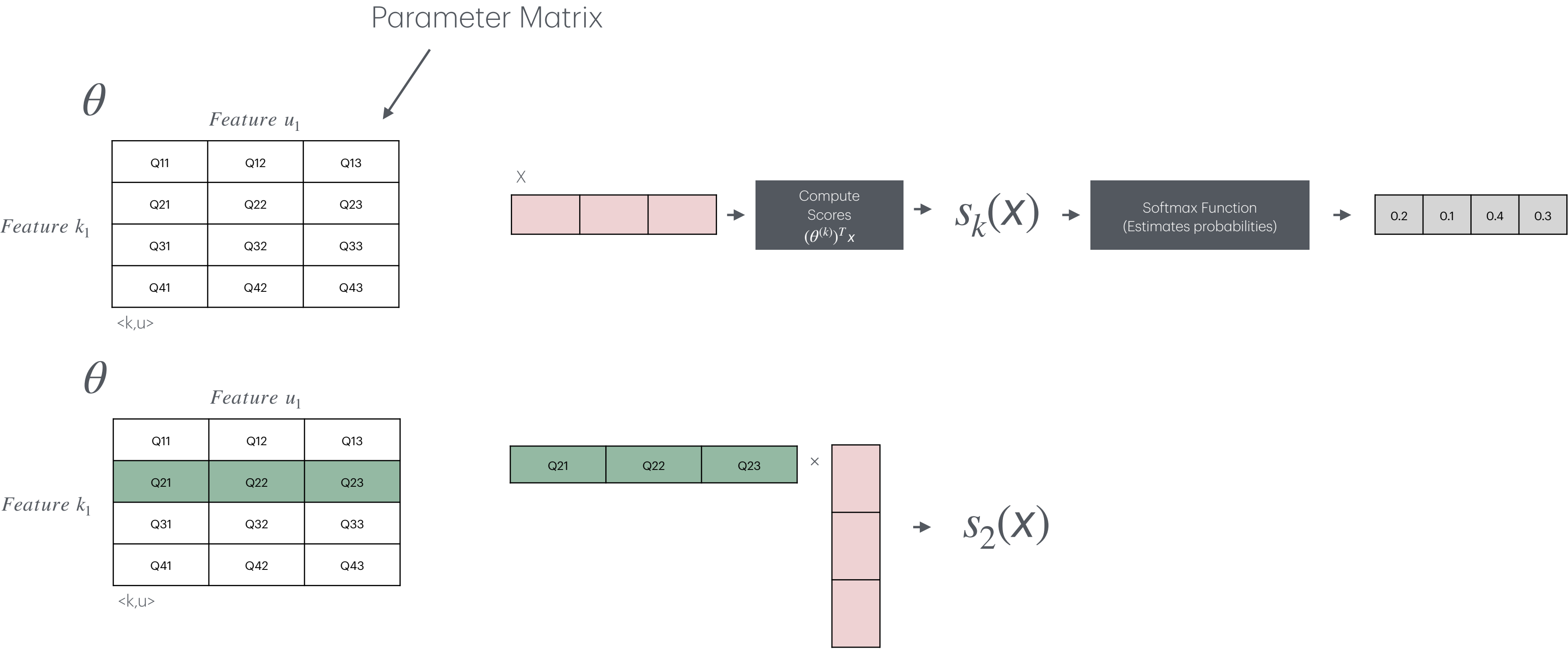
# Logistic Regression: How to Train?

No closed form for Logistic Regression. Find partial derivatives of cost function and update parameters at each step (Batch, stochastic, or mini-batch)

0   1   1   1

$y_i = 0$

$x_i$

· · ·

| Random Instance | → | Gradient Descent | → | Predictions | → | Decision Boundary |

$\theta_2$

Class = 1

Class = 0

$\theta_1$

Green- 70%

Pink- 60%

Red - 40%

The positive class would be the probability of instance being class=1 (e.g. prob of being male given features). All other classes are NOT male.

# Softmax Regression: How to Train?

Parameter Matrix

| | $u_1$ | $u_2$ | $u_3$ | $k_1$ | $k_2$ | $k_3$ | $k_4$ |
|---|---|---|---|---|---|---|---|
| | X1 | X2 | X3 | Y1 | Y2 | Y3 | Y4 |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

$\theta$

*Feature $u_1$*

| Q11 | Q12 | Q13 |
|---|---|---|
| Q21 | Q22 | Q23 |
| Q31 | Q32 | Q33 |
| Q41 | Q42 | Q43 |

*Feature $k_1$*

<k,u>

$\theta$

*Feature $u_1$*

| Q11 | Q12 | Q13 |
|---|---|---|
| Q21 | Q22 | Q23 |
| Q31 | Q32 | Q33 |
| Q41 | Q42 | Q43 |

*Feature $k_1$*

<k,u>

X

Compute Scores $(\theta^{(k)})^T x$

$s_k(X)$

Softmax Function (Estimates probabilities)

| 0.2 | 0.1 | 0.4 | 0.3 |

| Q21 | Q22 | Q23 | × |

$s_2(X)$

$\hat{P}_k$ Estimated Probability

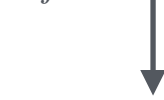| $s_1(X)$ | $s_2(X)$ | $s_3(X)$ | $s_4(X)$ |
|---|---|---|---|
| $\dfrac{exp(s_1(x))}{\sum_{j=1}^{K} exp(s_j(x))}$ | $\dfrac{exp(s_2(x))}{\sum_{j=1}^{K} exp(s_j(x))}$ | $\dfrac{exp(s_3(x))}{\sum_{j=1}^{K} exp(s_j(x))}$ | $\dfrac{exp(s_4(x))}{\sum_{j=1}^{K} exp(s_j(x))}$ |
| Estimated probability | Estimated probability | Estimated probability | Estimated probability |

Softmax predicts only one class

Parameter Matrix is updated during training, updating matrix during each training step that minimizes the cost function
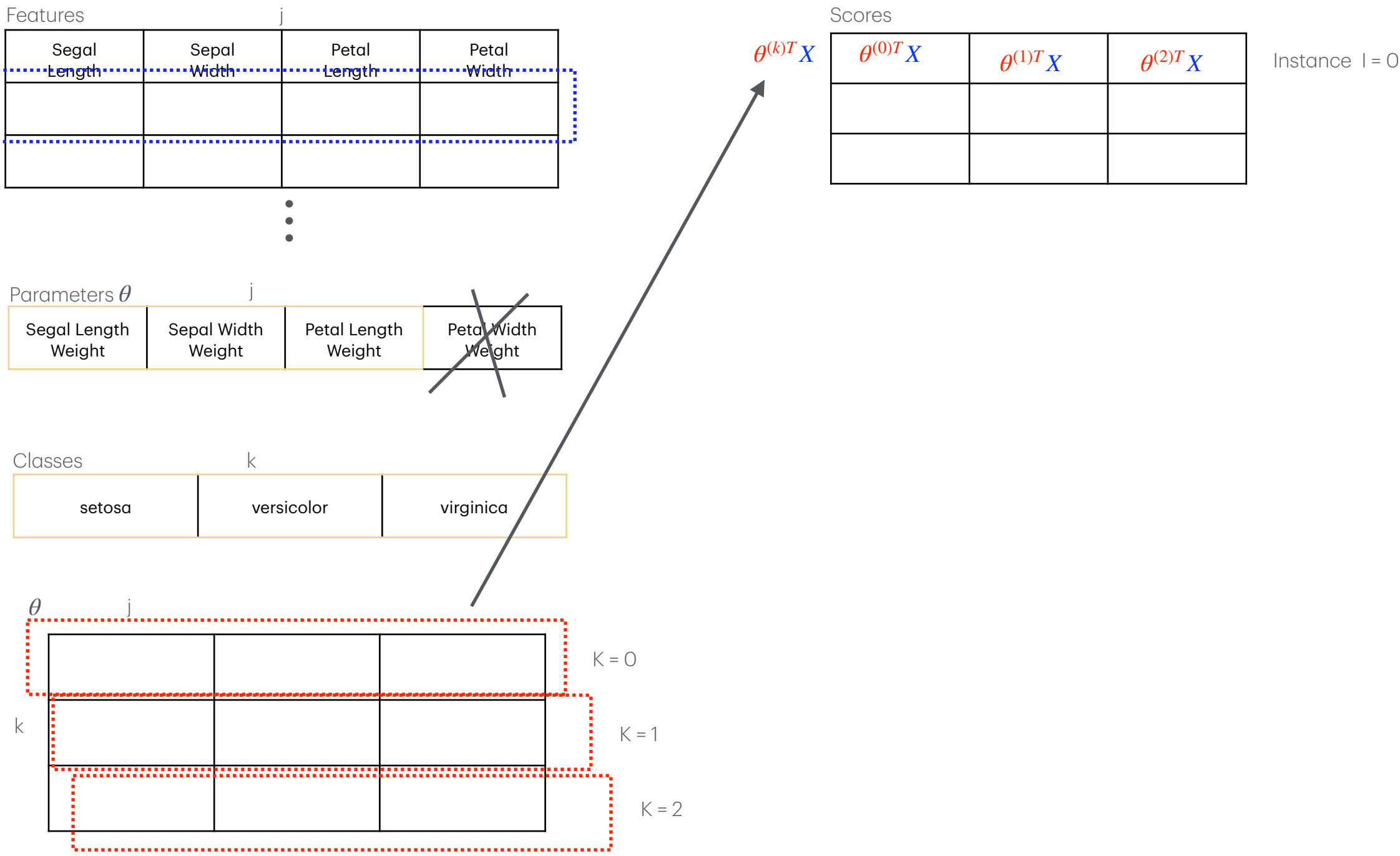
Prediction  $y' = argmax(estimated_{probabilities})$ ⟶ k index

# Batch Gradient Descent Example : Softmax

## Features

| | | j | |
|---|---|---|---|
| Segal Length | Sepal Width | Petal Length | Petal Width |
| | | | |
| | | | |

i

.
.
.

## Scores

| | | | |
|---|---|---|---|
| $\theta^{(k)T}X$ | $\theta^{(0)T}X$ | $\theta^{(1)T}X$ | $\theta^{(2)T}X$ |
| | | | |
| | | | |

Instance I = 0

## Parameters $\theta$

| | j | | |
|---|---|---|---|
| Segal Length Weight | Sepal Width Weight | Petal Length Weight | Petal Width Weight |

## Classes

| | k | |
|---|---|---|
| setosa | versicolor | virginica |

$\theta$

j

| | | |
|---|---|---|
| | | |
| | | |
| | | |

K = 0

k

K = 1

K = 2

$\theta$ · X = $class_{scores\_per\_instance}$

X · $\theta$ = $class_{scores\_all}$

$\theta$  ⌐  X

k

·

# Batch Gradient Descent Example : Softmax : Scores

Scores

exp(scores)

Exponential

Sum
Axis = 1

Numerator

Div

Denominator

$\hat{P}_k$

| 0.3 | 0 | 0.7 |
| 0 | 1 | 0 |
| 0.6 | 0.2 | 0.2 |

ArgMAX

$\hat{y}$  Predictions:: possible values are 0, 1, 2

Negative log used for cost because high probability will produce low value, and low probability produce high value

★ -log for these values are <= 0.5 may seem good, but we would seek ever lower values during gradient descent

- log

$-log(\hat{P}_k)$

| 1.2 | 6.9 | 0.3 ★ |
| 2.3 | 0.1 | 2.3 |
| 0.5 ★ | 1.69 | 1.609 |

$y$  Target

One hot expansion

$Cost = J(\theta) = CrossEntropy$

$y_{vectorize}$

| 0 | 0 | 1 ★ |
| 0 | 1 | 0 |
| ★ 1 | 0 | 0 |

X

Sum_all

$\dfrac{1}{num\_instances}$

Cost

# Batch Gradient Descent Example : Softmax : Scores : Regularization

- log

$-log(\hat{P_k})$

| 1.2 | 6.9 | 0.3 ★ |
|-----|-----|-------|
| 2.3 | 0.1 | 2.3 |
| 0.5 ★ | 1.69 | 1.609 |

$y$ Target

One hot expansion

$y_{vectorize}$

| 0 | 0 | 1 ★ |
|---|---|-----|
| 0 | 1 | 0 |
| 1 ★ | 0 | 0 |

X

k=1  k=2  k=3

$Sum_{axis=0}$

Regularization
hyperparameters

k=1

$\alpha \sum_{j=1}^{4} |\theta_{k=1_j}|$

k=2

$\alpha \sum_{j=1}^{4} |\theta_{k=2_j}|$

k=3

$\alpha \sum_{j=1}^{4} |\theta_{k=3_j}|$

+

Sum_all

$\dfrac{1}{num\_instances}$
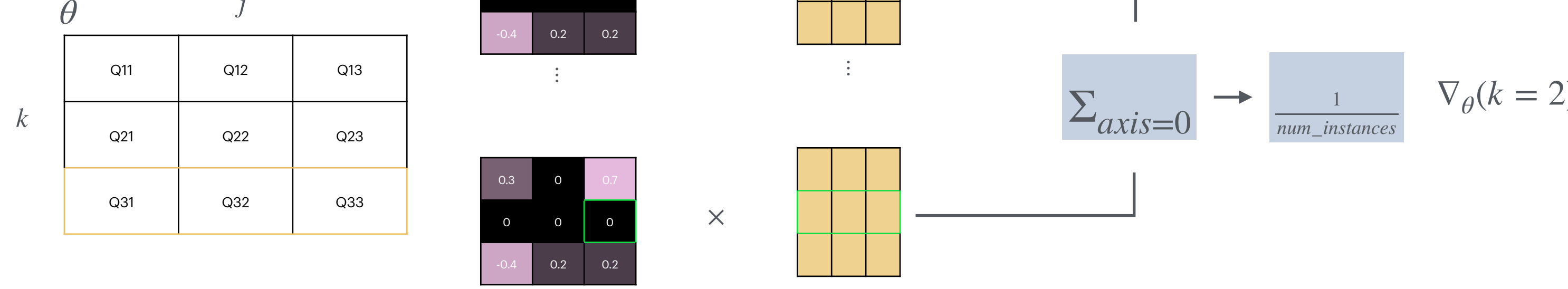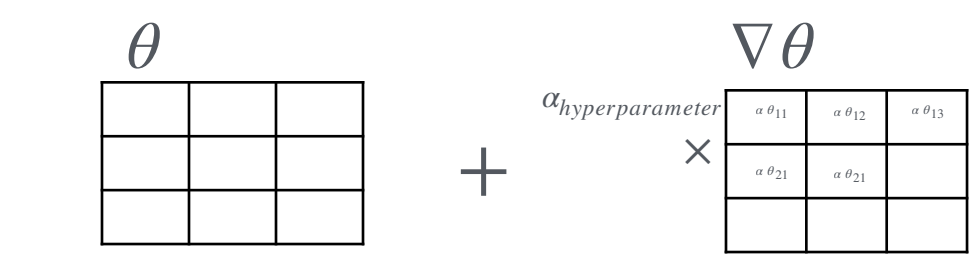
Cost

# Gradient  Vector Cost



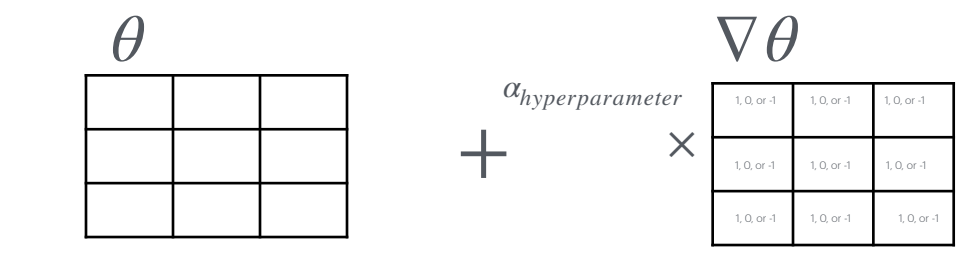Gradient vector for class k = 0

Gradient vector for class k = 1

Gradient vector for class k = 2

$\hat{P}_k - y_{vectorize}$

$X$

$\hat{P}_k$
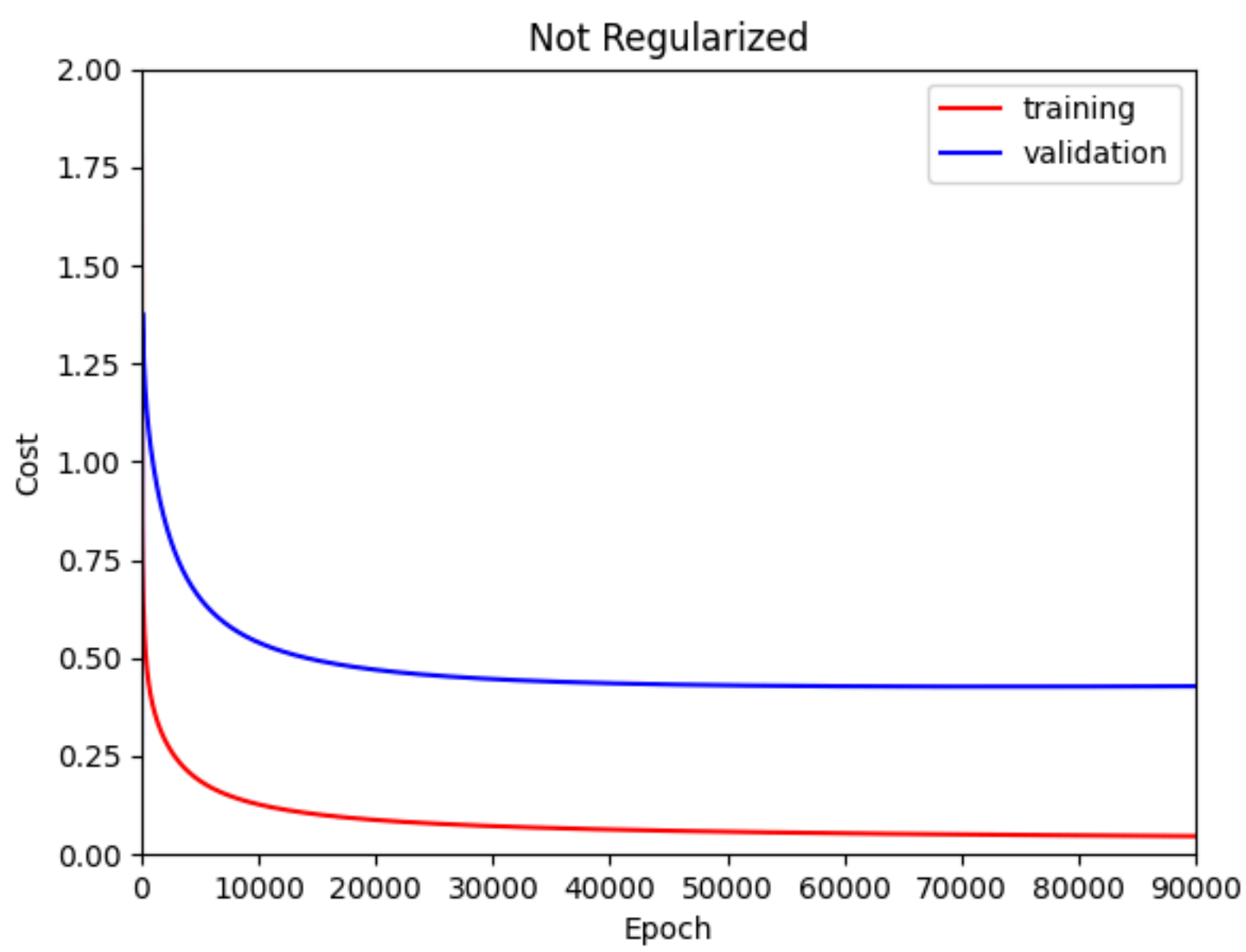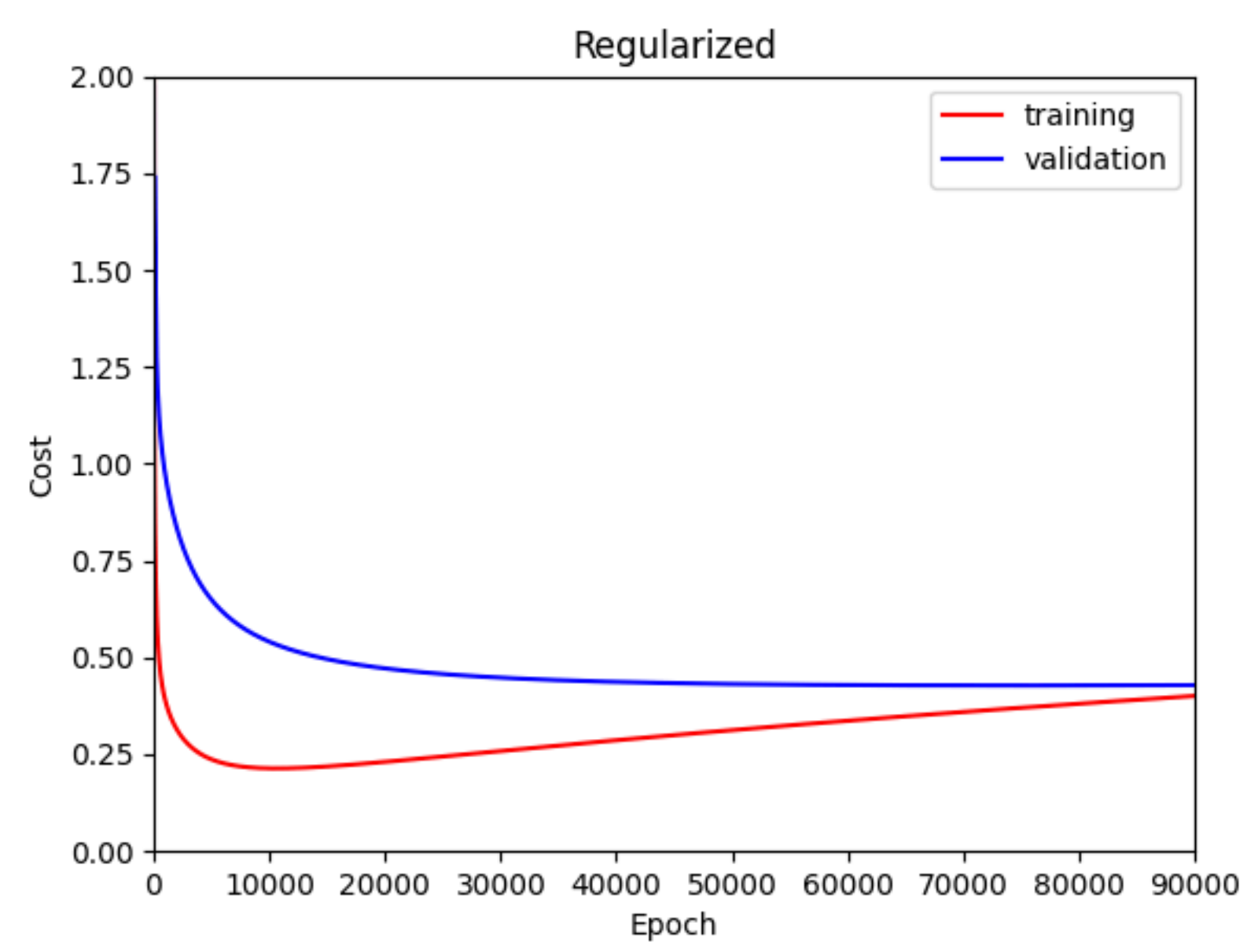
$y_{vectorize}$

Target

$y$

One hot expansion

$\hat{P}_k - y_{vectorize} = error_{matrix}$

$X$

$\theta$      $j$

$k$

| Q11 | Q12 | Q13 |
| Q21 | Q22 | Q23 |
| Q31 | Q32 | Q33 |

$\Sigma_{axis=0}$

$\frac{1}{num\_instances}$

$\nabla_\theta(k=0)$

$\nabla_\theta(k=1)$

$\nabla_\theta(k=2)$

Ridge Gradient update

$\theta$

$\nabla\theta$

$\alpha_{hyperparameter}$

Lasso Gradient update

$\theta$

$\nabla\theta$

$\alpha_{hyperparameter}$

Lasso

Ridge