# Linear models I
## Regression
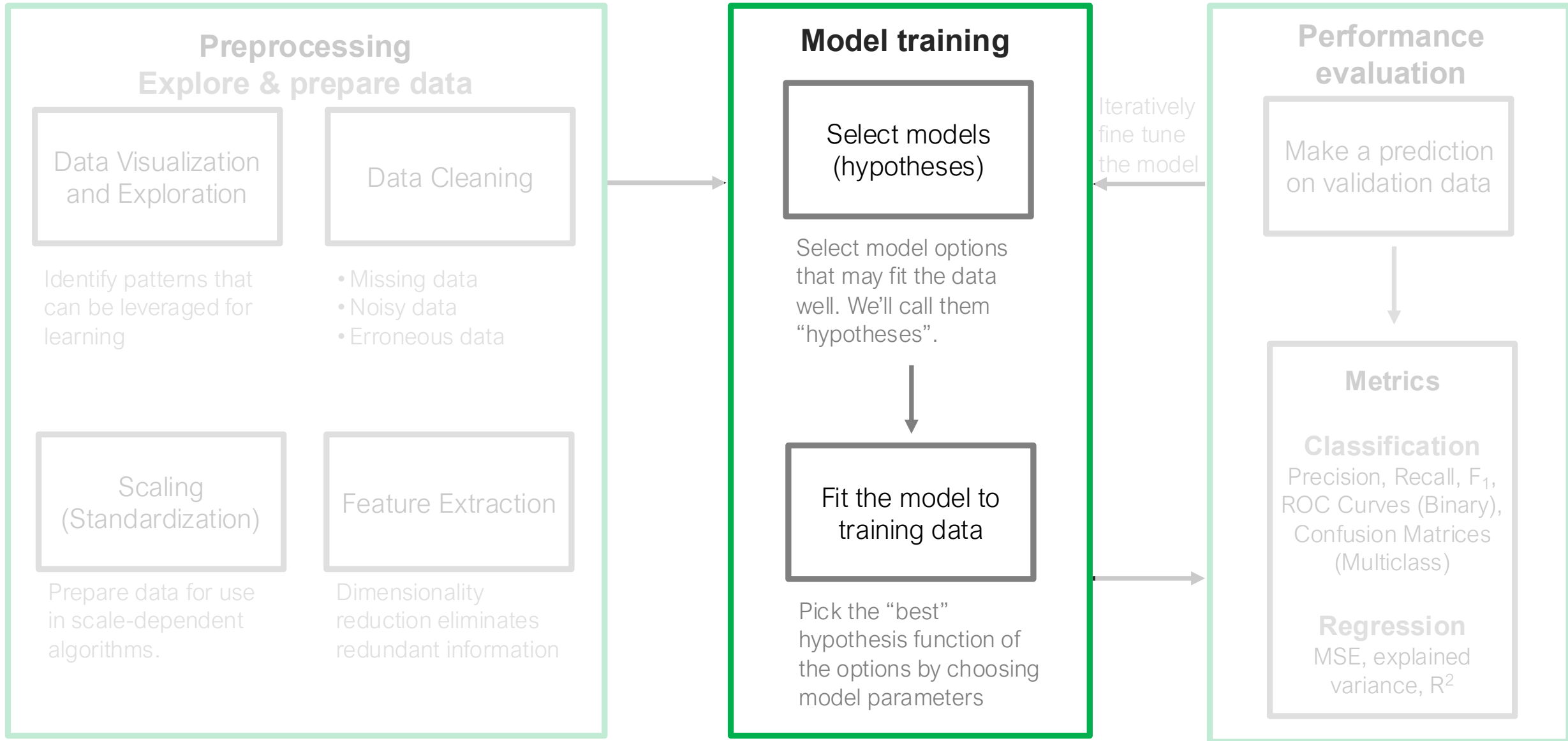
# Supervised learning in practice

## Preprocessing
**Explore & prepare data**

**Data Visualization and Exploration**

Identify patterns that can be leveraged for learning

**Data Cleaning**

• Missing data
• Noisy data
• Erroneous data

**Scaling (Standardization)**

Prepare data for use in scale-dependent algorithms.

**Feature Extraction**

Dimensionality reduction eliminates redundant information

## Model training

**Select models (hypotheses)**

Select model options that may fit the data well. We'll call them "hypotheses".

**Fit the model to training data**

Pick the "best" hypothesis function of the options by choosing model parameters

Iteratively fine tune the model

## Performance evaluation

**Make a prediction on validation data**

**Metrics**

**Classification**
Precision, Recall, $F_1$, ROC Curves (Binary), Confusion Matrices (Multiclass)

**Regression**
MSE, explained variance, $R^2$

# Supervised learning in practice

**Preprocessing**
**Explore & prepare data**

Data Visualization and Exploration

Data Cleaning

Identify patterns that can be leveraged for learning

- Missing data
- Noisy data
- Erroneous data

Scaling (Standardization)

Feature Extraction

Prepare data for use in scale-dependent algorithms.

Dimensionality reduction eliminates redundant information

**Model training**

**1** Supervised Learning Models: Linear models and KNN
(enough to get started using supervised learning)

**5** Other algorithms and concepts:
- Generative vs discriminative models
- Parametric vs nonparametric models
- Model ensembles
- Feature/representation learning (neural networks, deep learning)

**4** How to control model overfit: regularization strategies for model refinement

Iteratively fine tune the model

**Performance evaluation**

Make a prediction on validation data

**2** Evaluating model performance and comparing models

**Classification**

**3** How to make decisions using models

**Regression**
MSE, explained variance, $R^2$

# Supervised learning in practice

### Preprocessing
### Explore & prepare data

**Data Visualization and Exploration**

Identify patterns that can be leveraged for learning

**Data Cleaning**

• Missing data
• Noisy data
• Erroneous data

**Scaling (Standardization)**

Prepare data for use in scale-dependent algorithms.

**Feature Extraction**

Dimensionality reduction eliminates redundant information

### Model training

**Select models (hypotheses)**

Select model options that may fit the data well. We'll call them "hypotheses".

**Fit the model to training data**

Pick the "best" hypothesis function of the options by choosing model parameters

Iteratively fine tune the model

### Performance evaluation

Make a prediction on validation data

**Metrics**

**Classification**
Precision, Recall, $F_1$, ROC Curves (Binary), Confusion Matrices (Multiclass)

**Regression**
MSE, explained variance, $R^2$

# Model Fitting / Training Process

1. Choose a **hypothesis set of models** to train
   (e.g. linear regression with 4 predictor variables)

2. Identify a **cost function** to measure the model fit to the training data
   (e.g. mean square error)

3. **Optimize** model **parameters** to minimize cost
   (e.g. closed form solution using the normal equations for OLS)

★ We will use this procedure for ALL the parametric models we encounter
   Parametric models = models where the number of parameters are fixed and independent of the training data size

# Simple Linear Regression Model Fitting Process

$$\hat{y} = \phi_0 + \phi_1 x$$



Images from Prince, S.J., 2023. *Understanding Deep Learning*. MIT press.

a) Loss, $L[\phi]$

$$L[\phi] = \sum_{i=1}^{N} (\hat{y}_i - y_i)^2$$

$$\hat{y}_i = \phi_0 + \phi_1 x_i$$

Slope, $\phi_1$

Intercept, $\phi_0$

b) Output, $y$

Input, $x$

Images from Prince, S.J., 2023. *Understanding Deep Learning*. MIT press.

b) Loss, $L[\phi]$

Slope, $\phi_1$ vs Intercept, $\phi_0$

Loss, $L = 7.07$
$\phi_0 = 0.4, \phi_1 = 0.2$

c) Loss, $L = 10.28$
$\phi_0 = 1.60, \phi_1 = -0.8$

d) Loss, $L = 0.20$
$\phi_0 = 0.84, \phi_1 = 0.5$

Output, $y$ — Input, $x$

Images from Prince, S.J., 2023. *Understanding Deep Learning*. MIT press.

# How can we…

define what makes a model linear?

fit our model to our training data?

# What makes a model **linear**?

Linear models

# Which of the following models are linear?

Model parameters are $w_i$

Target variable is $y$

Remaining components are features

A $\quad y = w_0$

B $\quad y = w_0 + w_1 x_1$

C $\quad y = w_0 + w_1 x_1 + x_2^{w_2}$

D $\quad y = w_0 + w_1 x_1^2 + w_2 x_2^{0.4}$

E $\quad y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1 x_2$

F $\quad y = w_0 + w_1 \int \sqrt[3]{x_1}\, dx_1 + w_2 g(x_2) + w_3 median(x_1, x_2, x_3)$

# Linear models are linear in the parameters

A linear combination is quantity where a set of terms are added together, each multiplied by a constant (parameter)

$$\boldsymbol{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \qquad \boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

In other words, a linear combination is a sum of scalar multiples of vectors:

$$\boldsymbol{z} = 2\boldsymbol{x} + 99.9999\boldsymbol{y}$$

$$\boldsymbol{z} = 0.333\boldsymbol{x} + (-42)\boldsymbol{y}$$

$$\boldsymbol{z} = 0\boldsymbol{x} + \left(\frac{\pi^2}{e}\right)\boldsymbol{y}$$

# Linear regression model

Parameters $\quad w_j$

Target $\qquad y_i$

Features $\qquad x_{i,j}$

Sample $\qquad$ Feature

$$y_i = \sum_{j=0}^{p} w_j x_{i,j} = \mathbf{w}^\top \mathbf{x}_i$$

$$y_i = w_0 x_{i,0} + w_1 x_{i,1} + w_2 x_{i,2} + \cdots + w_p x_{i,p}$$

Intercept/bias term: $\quad x_{i,0} \triangleq 1$

$$y_i = w_0 + w_1 x_{i,1} + w_2 x_{i,2} + \cdots + w_p x_{i,p}$$

# Which of the following models are linear?

…or could be easily framed as…

Model parameters are $w_i$

Target variable is $y$

Remaining components are features

**A**  $y = w_0$

**B**  $y = w_0 + w_1 x_1$

**C**  $y = w_0 + w_1 x_1 + x_2^{w_2}$

**D**  $y = w_0 + w_1 x_1^2 + w_2 x_2^{0.4}$

**E**  $y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1 x_2$

**F**  $y = w_0 + w_1 \int \sqrt[3]{x_1}\, dx_1 + w_2 g(x_2) + w_3 median(x_1, x_2, x_3)$

# Which of the following models are linear?

A   $y = w_0$

B   $y = w_0 + w_1 x_1$

C   $y = w_0 + w_1 x_1 + x_2^{w_2}$

D   $y = w_0 + w_1 x_1^2 + w_2 x_2^{0.4}$

E   $y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1 x_2$

F   $y = w_0 + w_1 \int \sqrt[3]{x_1} \, dx_1 + w_2 g(x_2) + w_3 median(x_1, x_2, x_3)$

…or could be easily framed as…

Model parameters are $w_i$

Target variable is $y$

Remaining components are features

**ALL except C** are linear
in the **parameters, $w$**

# We can rewrite these models in terms of transformed features

It becomes clear this is a linear model when we rewrite the features:

$$y = w_0 + w_1 z_1 + w_2 z_2 + w_3 z_3$$

Transformed features:

$$z_1 = \int \sqrt[3]{x_1} dx_1 \qquad z_2 = g(x_2) \qquad z_2 = median(x_1, x_2, x_3)$$

F $$y = w_0 + w_1 \int \sqrt[3]{x_1} dx_1 + w_2 g(x_2) + w_3 median(x_1, x_2, x_3)$$

# Notation

$x_{i,j}$

Sample    Feature

**Number of features**

**Number of samples**

| | 1 | | $p$ | |
|---|---|---|---|---|
| | | Shorthand | | Shorthand — Assume $x_{i,0} = 1$ |
| 1 | $x_{i,j}$ | $x_i$ | $\boldsymbol{x}_i = \begin{bmatrix} x_{i,0} \\ x_{i,1} \\ \vdots \\ x_{i,p} \end{bmatrix}$ | $\boldsymbol{x}_i = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_p \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{bmatrix}$ |
| $N$ | $\boldsymbol{x} = \begin{bmatrix} x_{1,j} \\ x_{2,j} \\ \vdots \\ x_{N,j} \end{bmatrix}$ | $\boldsymbol{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$ | $\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1^T \\ \boldsymbol{x}_2^T \\ \vdots \\ \boldsymbol{x}_N^T \end{bmatrix} = \begin{bmatrix} x_{1,0} & x_{1,1} & \dots & x_{1,p} \\ x_{2,0} & x_{2,1} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,0} & x_{N,1} & \dots & x_{N,p} \end{bmatrix}$ | |

# Types of Linear Regression

Here, $x_0 \triangleq 1$

| | **One** feature variable | **Two or more** feature variables |
|---|---|---|
| **One** target variable | **Simple Linear Regression** $y_i = w_0 + w_1 x_{i,1}$ | **Multiple Linear Regression** $y_i = \sum_{j=0}^{p} w_j x_{i,j}$ or $y_i = \boldsymbol{w}^\top \boldsymbol{x}_i$ |
| **Two or more** target variables | **Multivariate (Multiple) Linear Regression** | |

q = # target variables
[q x 1]

$$\boldsymbol{y}_i = \boldsymbol{W}^\top \boldsymbol{x}_i$$

[q x p]  [p x 1]

p = # features (includes $x_0 = 1$)

# Linear models: pros and cons

**Pros**

Simple/fast to implement and interpret

Excels if the relationship between features and targets is linear or can be expressed in terms of linear combinations of features

Often a good starting point or baseline model for many analyses

**Cons**

For many applications with complex feature-target relationships, underfits

Requires feature engineering for capturing more complex feature-target relationships

# How do we fit the model to the training data?

A winding path to the least squares solution

# Linear models and error

$\hat{f}(x)$

$y$

$E(x_i) = \hat{f}(x_i) - y_i$

$x_i$

$x$

**simple linear regression**

$\hat{f}(\boldsymbol{x})$

$y$

$x_1$

$x_2$

**multiple linear regression**

# How do we fit a linear model to training data?

Training data:
$$D = \{(\boldsymbol{x}_1, y_1), (\boldsymbol{x}_2, y_2), \dots, (\boldsymbol{x}_N, y_N)\}$$

We want the error between our predictions and training data to be small



$\hat{f}(x)$

$y$

Predictions

Training labels

$$E_i = \hat{f}(x_i) - y_i$$

**Minimize this!**

$x_i$

$x$

# How do we measure error?

The error for one sample is:

$$E_i = \hat{f}(\boldsymbol{x}_i) - y_i$$

$$\hat{y}_i = \hat{f}(\boldsymbol{x}_i) = \sum_{j=0}^{p} w_j x_{i,j}$$

We want to minimize the error across all our training data, so…
we use mean squared error to quantify training (in-sample) error:

**Training (in-sample) error:** $\quad E_{in}(\hat{f}, D) = \dfrac{1}{N} \sum_{n=1}^{N} \left( \hat{f}(\boldsymbol{x}_n) - y_n \right)^2$

We call this our **Cost Function**  (a.k.a. loss, error, or objective)

Why mean squared error? See *Understanding Deep Learning* Chapter 5 (particularly 5.3)

**Cost Function:**

$$C(\hat{f}, D) = E_{in}(\hat{f}, D) = \frac{1}{N}\sum_{n=1}^{N}\left(\hat{f}(x_n) - y_n\right)^2$$

where $\hat{f}$ is the **model** and $D$ is the **training data**.

Training error is a function of our **model** and the **training data**

We can't change the data; we adjust our model to minimize cost

To adjust the model, we choose model **parameters** that minimize cost

This is an **optimization** problem

# How to fit our model to the training data?

Equivalently: how do we choose $\boldsymbol{w}$ to minimize cost (error)

$$E_{in}(\hat{f}, D) = \frac{1}{N} \sum_{n=1}^{N} \left(\hat{f}(\boldsymbol{x}_n) - y_n\right)^2$$

where $\hat{f}(\boldsymbol{x}_n) = \boldsymbol{w}^T \boldsymbol{x}_n$

We want to minimize

...by varying $\boldsymbol{w}$

How do we do that?

$$E_{in}(\boldsymbol{w}) = \frac{1}{N} \sum_{n=1}^{N} (\boldsymbol{w}^T \boldsymbol{x}_n - y_n)^2$$

Once we've determined the form of the model and the training data, the model parameters $\boldsymbol{w}$ are the only parts we can adjust

# Calculus

# A moment of calculus

Function of one variable

$$f(x) = ax + bx^2$$

Function of multiple variables

$$f(x_1, x_2) = ax_1 + bx_2$$

| **Derivative** | **Partial Derivative** | **Gradient** |
|:---:|:---:|:---:|

$$\frac{df}{dx} = a + 2bx$$

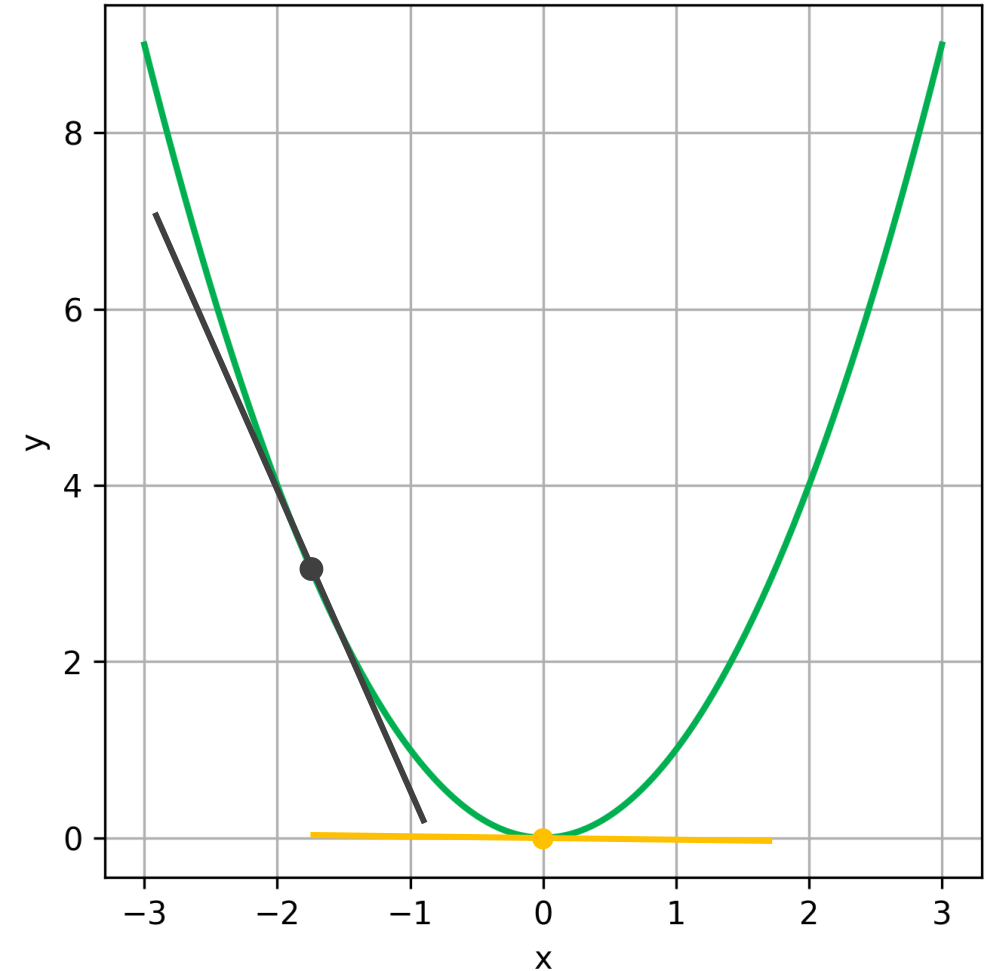$$\frac{\partial f}{\partial x_1} = a$$

$$\frac{\partial f}{\partial x_2} = b$$

$$\nabla_x f = \begin{bmatrix} \dfrac{\partial f}{\partial x_1} \\ \dfrac{\partial f}{\partial x_2} \end{bmatrix}$$

$$= \begin{bmatrix} a \\ b \end{bmatrix}$$

May also treat parameters as variables and take their partial derivative   $\frac{\partial f}{\partial b} = x_2$

# Why set the derivative to 0?

Remember the derivative (and gradient) is the direction and rate of fastest increase

For a continuous, convex function, the function is minimized when that derivative is zero

# How to fit our model to the training data?

Take the gradient with respect to $\boldsymbol{w}$, set it to zero, and solve for $\boldsymbol{w}$
(think derivative)

$$\nabla_{\boldsymbol{w}} E_{in}(\boldsymbol{w}) = \nabla_{\boldsymbol{w}} \left( \frac{1}{N} \sum_{n=1}^{N} (\boldsymbol{w}^T \boldsymbol{x}_n - y_n)^2 \right)$$

$p$ = number of predictors
$N$ = number of data points

We will walk through the **ordinary least squares** (OLS) closed-form solution.

$$\nabla_{\boldsymbol{w}} E_{in}(\boldsymbol{w}) = \begin{bmatrix} \dfrac{\partial E_{in}}{\partial w_0} \\ \dfrac{\partial E_{in}}{\partial w_1} \\ \vdots \\ \dfrac{\partial E_{in}}{\partial w_p} \end{bmatrix} = \boldsymbol{0}$$

Size: $[p + 1 \times 1]$ or $\mathbb{R}^{p+1\times1}$

We could have used another optimization approach like **gradient descent**

# How to fit our model to the training data?

Our cost function…

$$E_{in}(\boldsymbol{w}) = \frac{1}{N} \sum_{n=1}^{N} (\boldsymbol{w}^T \boldsymbol{x}_n - y_n)^2$$

Scalar

$$\boldsymbol{w}^T \in \mathbb{R}^{1 \times p+1}$$
$$\boldsymbol{x}_n \in \mathbb{R}^{p+1 \times 1}$$

…can be rewritten as:

$$E_{in}(\boldsymbol{w}) = \frac{1}{N} (\boldsymbol{Xw} - \boldsymbol{y})^T (\boldsymbol{Xw} - \boldsymbol{y})$$

Convenient definitions:

$$\boldsymbol{y} \in \mathbb{R}^{N \times 1}$$

$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1^T \\ \boldsymbol{x}_2^T \\ \vdots \\ \boldsymbol{x}_N^T \end{bmatrix} \in \mathbb{R}^{N \times p+1}$$

$$\boldsymbol{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_p \end{bmatrix} \in \mathbb{R}^{p+1 \times 1}$$

$p$ = number of predictors
$N$ = number of data points

**①** Assume $p = 2$, $N = 4$

$p$ = number of predictors
$N$ = number of data points

$$\boldsymbol{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \qquad \boldsymbol{x}_i = \begin{bmatrix} x_{i,0} \\ x_{i,1} \\ x_{i,2} \end{bmatrix}$$

**②**

$$\boldsymbol{w}^T \boldsymbol{x}_i = \begin{bmatrix} w_0 & w_1 & w_2 \end{bmatrix} \begin{bmatrix} x_{i,0} \\ x_{i,1} \\ x_{i,2} \end{bmatrix}$$

$$= w_0 x_{i,0} + w_1 x_{i,1} + w_2 x_{i,2}$$

**③**

$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1^T \\ \boldsymbol{x}_2^T \\ \boldsymbol{x}_3^T \\ \boldsymbol{x}_4^T \end{bmatrix} = \begin{bmatrix} x_{1,0} & x_{1,1} & x_{1,2} \\ x_{2,0} & x_{2,1} & x_{2,2} \\ x_{3,0} & x_{3,1} & x_{3,2} \\ x_{4,0} & x_{4,1} & x_{4,2} \end{bmatrix}$$

**Aside on algebraic manipulations**

**④**

$$\boldsymbol{X}\boldsymbol{w} = \begin{bmatrix} x_{1,0} & x_{1,1} & x_{1,2} \\ x_{2,0} & x_{2,1} & x_{2,2} \\ x_{3,0} & x_{3,1} & x_{3,2} \\ x_{4,0} & x_{4,1} & x_{4,2} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} \boldsymbol{w}^T \boldsymbol{x}_1 \\ \boldsymbol{w}^T \boldsymbol{x}_2 \\ \boldsymbol{w}^T \boldsymbol{x}_3 \\ \boldsymbol{w}^T \boldsymbol{x}_4 \end{bmatrix}$$

**4**

$$\boldsymbol{Xw} = \begin{bmatrix} x_{1,0} & x_{1,1} & x_{1,2} \\ x_{2,0} & x_{2,1} & x_{2,2} \\ x_{3,0} & x_{3,1} & x_{3,2} \\ x_{4,0} & x_{4,1} & x_{4,2} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} \boldsymbol{w}^T \boldsymbol{x}_1 \\ \boldsymbol{w}^T \boldsymbol{x}_2 \\ \boldsymbol{w}^T \boldsymbol{x}_3 \\ \boldsymbol{w}^T \boldsymbol{x}_4 \end{bmatrix}$$

**5**

$$\boldsymbol{Xw} - \boldsymbol{y} = \begin{bmatrix} \boldsymbol{w}^T \boldsymbol{x}_1 - \boldsymbol{y}_1 \\ \boldsymbol{w}^T \boldsymbol{x}_2 - \boldsymbol{y}_2 \\ \boldsymbol{w}^T \boldsymbol{x}_3 - \boldsymbol{y}_3 \\ \boldsymbol{w}^T \boldsymbol{x}_4 - \boldsymbol{y}_4 \end{bmatrix}$$

**6**

$$(\boldsymbol{Xw} - \boldsymbol{y})^T(\boldsymbol{Xw} - \boldsymbol{y}) = \begin{bmatrix} \boldsymbol{w}^T \boldsymbol{x}_1 - \boldsymbol{y}_1 & \boldsymbol{w}^T \boldsymbol{x}_2 - \boldsymbol{y}_2 & \boldsymbol{w}^T \boldsymbol{x}_3 - \boldsymbol{y}_3 & \boldsymbol{w}^T \boldsymbol{x}_4 - \boldsymbol{y}_4 \end{bmatrix} \begin{bmatrix} \boldsymbol{w}^T \boldsymbol{x}_1 - \boldsymbol{y}_1 \\ \boldsymbol{w}^T \boldsymbol{x}_2 - \boldsymbol{y}_2 \\ \boldsymbol{w}^T \boldsymbol{x}_3 - \boldsymbol{y}_3 \\ \boldsymbol{w}^T \boldsymbol{x}_4 - \boldsymbol{y}_4 \end{bmatrix}$$

$$= \sum_{n=1}^{N} (\boldsymbol{w}^T \boldsymbol{x}_n - \boldsymbol{y}_n)^2$$

**Aside on algebraic manipulations**

# How to fit our model to the training data?

Our cost function…

$$E_{in}(\boldsymbol{w}) = \frac{1}{N} \sum_{n=1}^{N} (\boldsymbol{w}^T \boldsymbol{x}_n - y_n)^2$$

↑
Scalar

$$\boldsymbol{w}^T \in \mathbb{R}^{1 \times p+1}$$
$$\boldsymbol{x}_n \in \mathbb{R}^{p+1 \times 1}$$

…can be rewritten as:

$$E_{in}(\boldsymbol{w}) = \frac{1}{N} (\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y})^T (\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y})$$

Convenient definitions:

$$\boldsymbol{y} \in \mathbb{R}^{N \times 1}$$

$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1^T \\ \boldsymbol{x}_2^T \\ \vdots \\ \boldsymbol{x}_N^T \end{bmatrix} \in \mathbb{R}^{N \times p+1}$$

$$\boldsymbol{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_p \end{bmatrix} \in \mathbb{R}^{p+1 \times 1}$$

$p$ = number of predictors
$N$ = number of data points

# How to fit our model to the training data?

$$E_{in}(\boldsymbol{w}) = \frac{1}{N}(X\boldsymbol{w} - \boldsymbol{y})^T(X\boldsymbol{w} - \boldsymbol{y})$$

(take gradient, set to 0)

$$\nabla_{\boldsymbol{w}} E_{in}(\boldsymbol{w}) = \frac{2}{N}(X^T X\boldsymbol{w} - X^T\boldsymbol{y}) = \boldsymbol{0}$$

(solve for $\boldsymbol{w}$)

**Univariate analogy**:

$$f(w) = \frac{1}{N}(xw - y)^2$$

$$= \frac{1}{N}(x^2 w^2 - 2xyw + y^2)$$

$$\frac{df(w)}{dw} = \frac{2}{N}(x^2 w - xy)$$

$$X^T X\boldsymbol{w} - X^T\boldsymbol{y} = \boldsymbol{0}$$

$$X^T X\boldsymbol{w} = X^T\boldsymbol{y} \quad \text{(normal equation)}$$

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

**Pseudoinverse** $\quad \mathrm{X}^\dagger = (X^T X)^{-1} X^T$

$$\boxed{\boldsymbol{w}^* = \mathrm{X}^\dagger \boldsymbol{y}}$$

# What is the pseudoinverse?

# Samples and features impact solutions

Features

Samples

$$\begin{bmatrix} N \times p \end{bmatrix} \begin{bmatrix} p \times 1 \end{bmatrix} = \begin{bmatrix} N \times 1 \end{bmatrix}$$

$$\boldsymbol{X} \qquad \boldsymbol{w} = \boldsymbol{y}$$

If $N = p$, then there are the same number of features as samples
(# equations = # unknowns)

If $N > p$, then the system of equations is **overdetermined**: more samples than features
(# equations > # unknowns)

Can't invert $\boldsymbol{X}$ – it's not square!

If $N < p$, then the system of equations is **underdetermined**: fewer samples than features
(# equations < # unknowns)

# Overdetermined systems

## Example 1

Fully determined system

# equations = # unknowns

$$w_1 = 2$$
$$w_2 = 1$$

Overdetermined system

# equations > # unknowns

$$w_1 = 2$$
$$w_2 = 1$$
$$w_2 = 5$$

## Example 2

Fully determined system

# equations = # unknowns

$$2w_1 + 3w_2 = 2$$
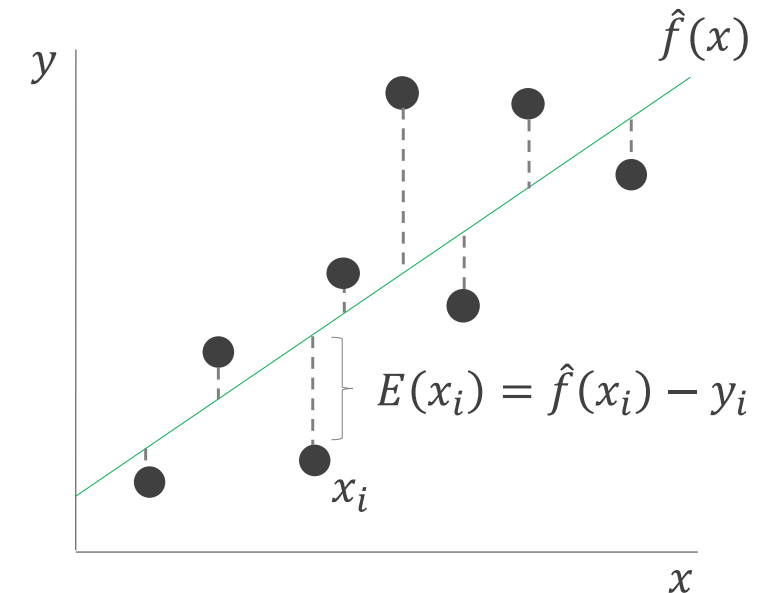$$6w_1 + 1w_2 = 0.63$$

Overdetermined system

# equations > # unknowns

$$2w_1 + 3w_2 = 2$$
$$6w_1 + 1w_2 = 0.63$$
$$3w_1 + 2w_2 = 14$$
$$16w_1 - w_2 = 0.1$$

$$E(x_i) = \hat{f}(x_i) - y_i$$

# Underdetermined systems

## Example 1

Fully determined system

# equations = # unknowns

$$w_1 = 2$$
$$w_2 = 1$$

Underdetermined system

# equations < # unknowns

$$w_1 = w_2$$

## Example 2

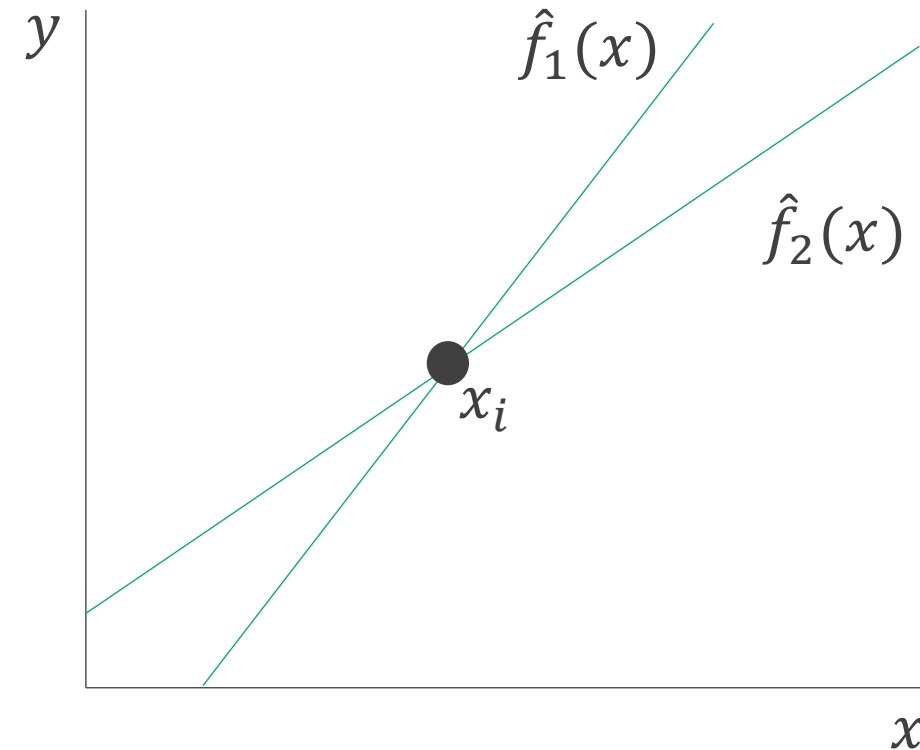Fully determined system

# equations = # unknowns

$$2w_1 + 3w_2 = 2$$
$$6w_1 + 1w_2 = 0.63$$

Underdetermined system

# equations < # unknowns

$$2w_1 + 3w_2 = 2$$

# What is the pseudoinverse?

Consider the case when $N = 3, p = 2$
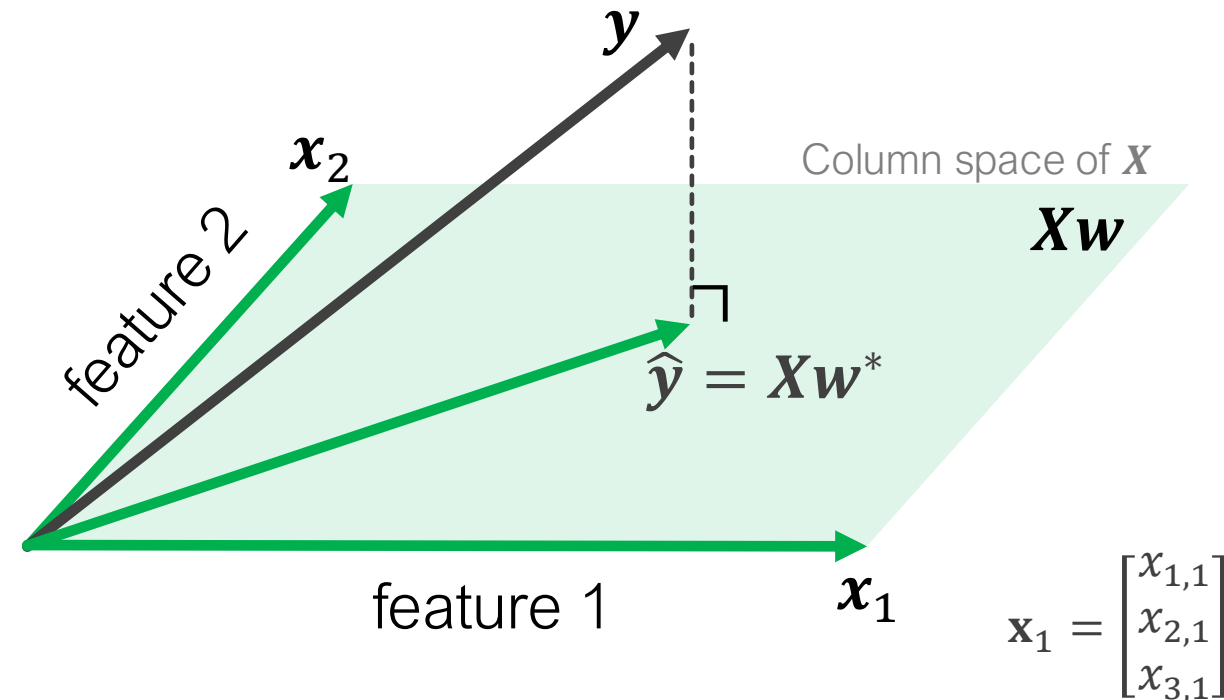(Assume no bias term here)

The least squares solution is the best we can do given $N > p$

Features

Samples
$$\begin{bmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \\ x_{3,1} & x_{3,2} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

Feature vectors: $x_1$      $x_2$

$$X \qquad w \quad \neq \quad y$$

We CAN use $X^\dagger$ to obtain $w^*$, the least squares solution:

$$w^* = X^\dagger y$$

$y$

$x_2$

Column space of $X$

$Xw$

feature 2

$\hat{y} = Xw^*$

feature 1                    $x_1$

$$\mathbf{x}_1 = \begin{bmatrix} x_{1,1} \\ x_{2,1} \\ x_{3,1} \end{bmatrix}$$

# Much of machine learning is optimizing a cost function

Least squares is one approach (when applicable), gradient descent is another, more generic method

# Model Fitting / Training Process

1.  Choose a **hypothesis set of models** to train
    (e.g. linear regression with 4 predictor variables)

2.  Identify a **cost function** to measure the model fit to the training data
    (e.g. mean square error)

3.  **Optimize** model **parameters** to minimize cost
    (e.g. closed form solution using the normal equations for OLS)

# We now have our model parameters, we can make predictions on unseen test data!

The parameters learned from model fitting

$$\hat{y}_i = \hat{f}(\boldsymbol{x}_i) = \sum_{j=0}^{p} w_j^* x_{i,j}$$

# Takeaways

Linear models are **linear in the weights**

Model fitting/training process (valid beyond linear models):
- Choose a hypothesis set of models to train
- Identify a cost function
- **Optimize the cost function** by adjusting model parameters
(This is the "learning" process)

Optimize cost functions for linear regression using least squares
- Least squares allows us to generate approximate solutions to overdetermined systems (more samples than features/parameters), which are common
- Alternative optimization strategies exist such as gradient descent