# Linear regression

Introduction to Machine Learning, Lecture 2

Fanny Yang

**D**INFK

**ETH** *zürich*

# Announcements

- Deadline for buddy matching: this **Friday 25.2. 23:59**

  - For today, the virtual attendees may take the time to brainstorm on their own

- Q&A session online **today 17:15 to 18:00**

  - may ask more administrative questions there

  - instructions on how to effectively use Moodle

  - find the link on the course website
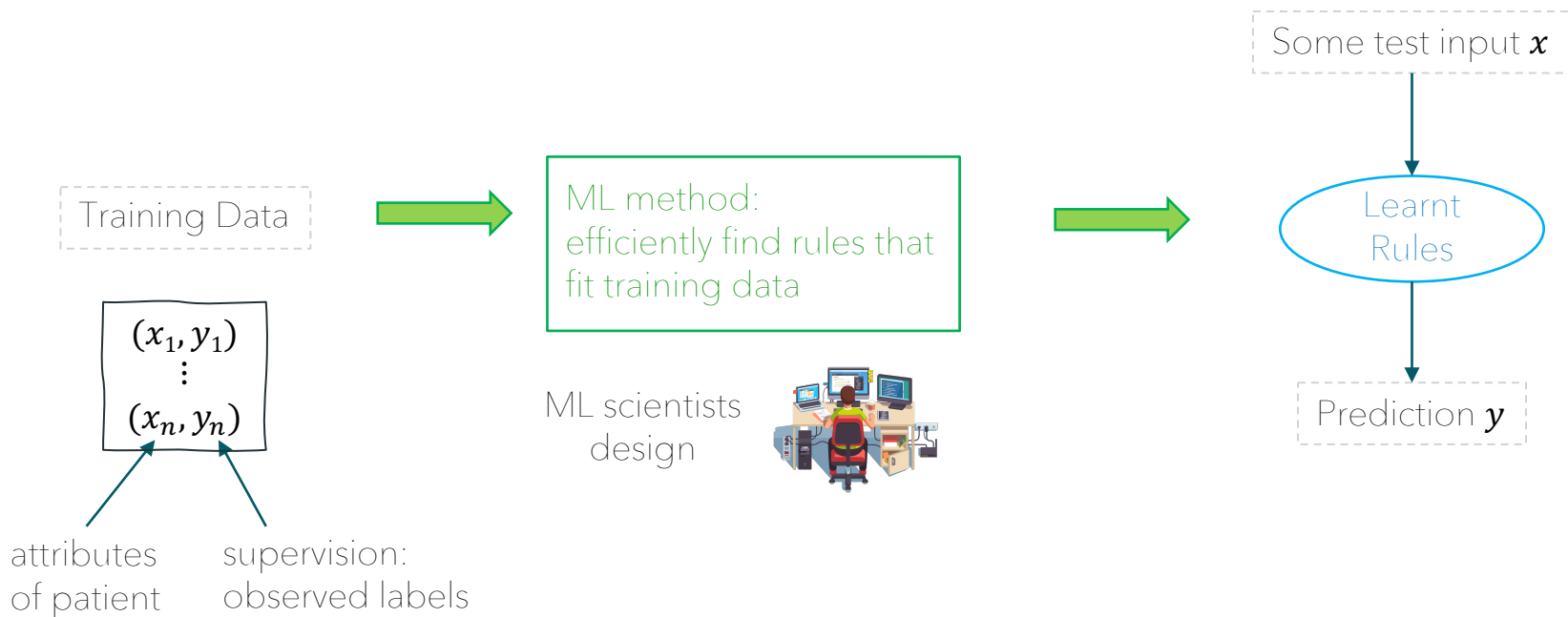
# Recap and plan for today

Last lecture: Different problems in machine learning
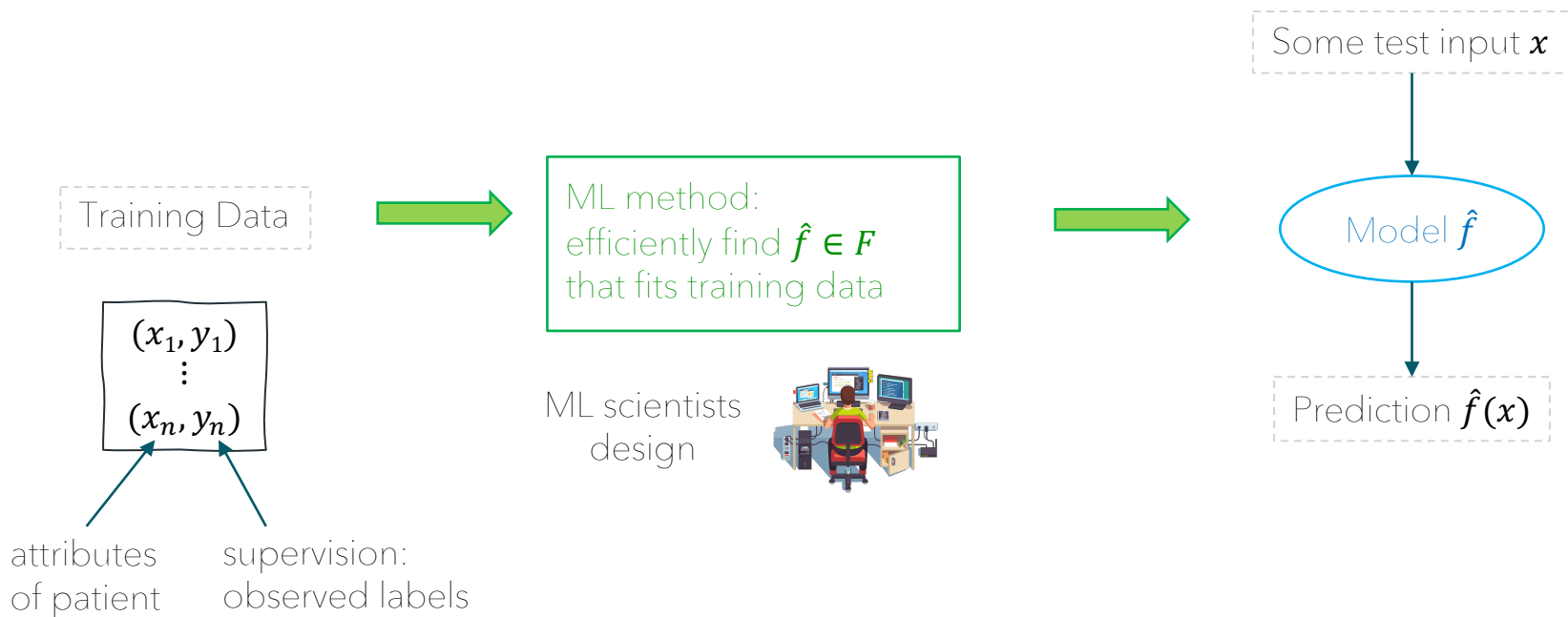
- supervised learning

- unsupervised learning

Today: Walking through the supervised learning pipeline with

- the simplest ML method: linear regression in 1d

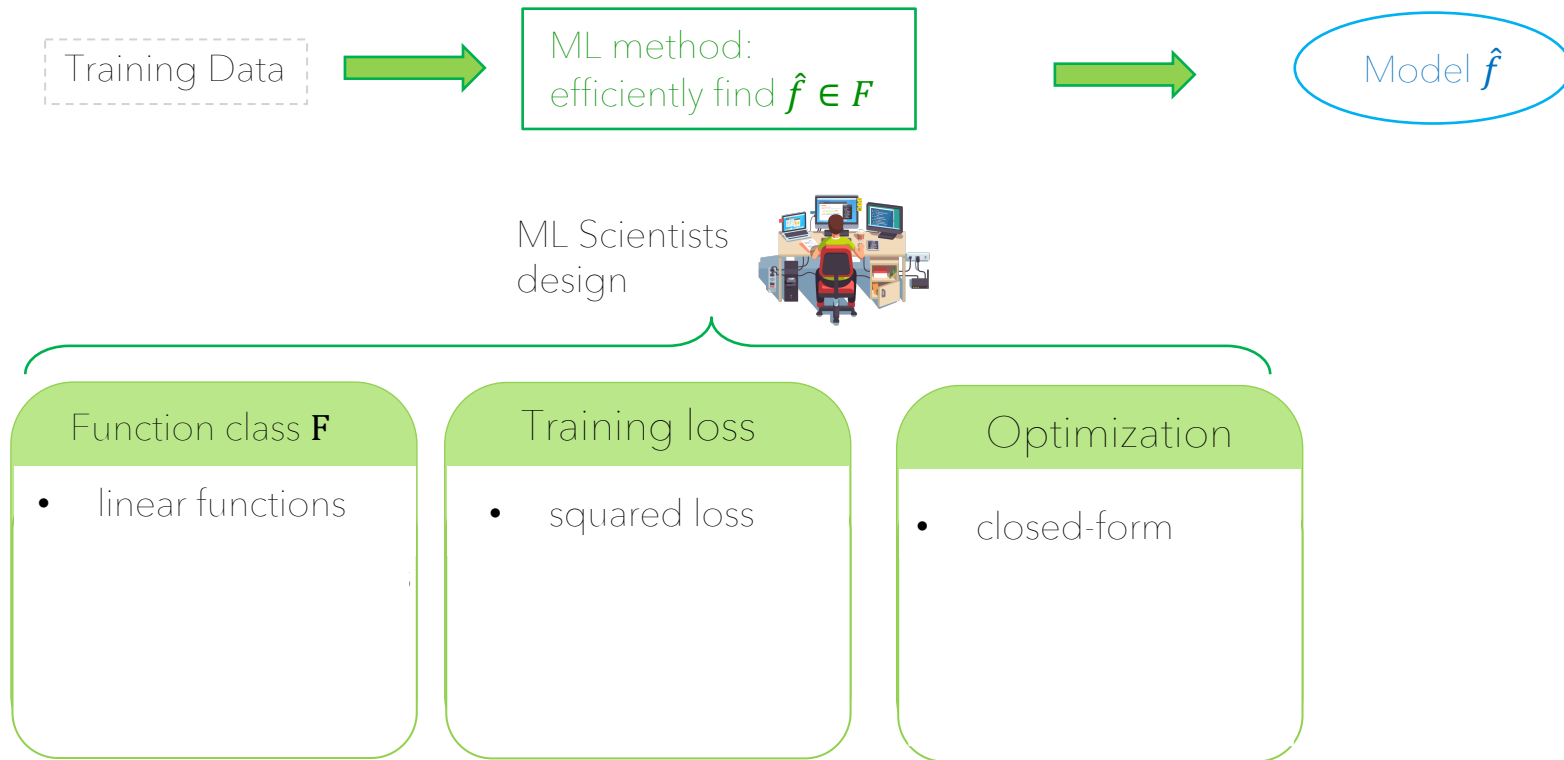- for the scenario of determining sales price of a house

# Simplified diagram of supervised learning

Training Data

$(x_1, y_1)$
$\vdots$
$(x_n, y_n)$

attributes
of patient

supervision:
observed labels

ML method:
efficiently find rules that
fit training data

ML scientists
design

Some test input $x$

Learnt
Rules

Prediction $y$

# Simplified diagram of supervised learning

Training Data

$(x_1, y_1)$
$\vdots$
$(x_n, y_n)$

attributes
of patient

supervision:
observed labels

ML method:
efficiently find $\hat{f} \in F$
that fits training data

ML scientists
design

Some test input $x$

Model $\hat{f}$

Prediction $\hat{f}(x)$

# Machine learning pipeline

Training Data → ML method: efficiently find $\hat{f} \in F$ → Model $\hat{f}$

ML Scientists design

| Function class $\mathbf{F}$ | Training loss | Optimization |
|---|---|---|
| • linear functions | • squared loss | • closed-form |

# Example (house)

# ML pipeline to find average house price

Last lecture I learned: I can use machine learning to do that. If I can find input–output examples, I can use supervised learning, it's a regression task!

Can't feed houses into computer! ⟶ Your house

Training Data ⟹

ML method:
efficiently find $\hat{f} \in F$
that fits training data

⟹ Model $\hat{f}$

Average price for your house in CHF

# **Step 0**: find representation for house

- Determine how to represent houses in a digital fashion

- for example using a vector of attributes:
  e.g. size, # bathrooms, distance to public transport, years since construction, …

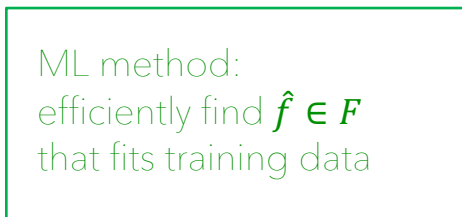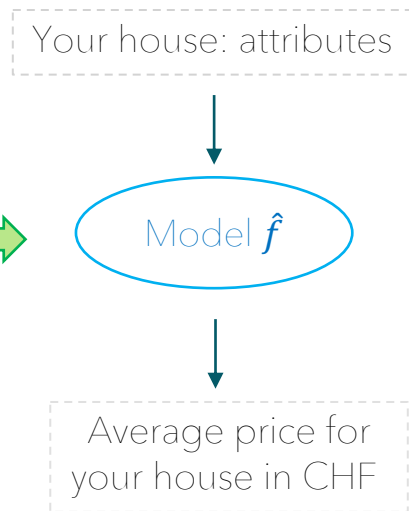- Collect attributes and sales prices from other houses and your own

extract →

*House 1: attributes, price*
*…*
*House n: attributes, price*

Your house: attributes

# ML pipeline to find average house price

**Step I: Collect**　　　　**Step II: Learn**　　　　**Step III: Predict**

Training Data

House 1: attributes, price
…
House n: attributes, price

ML method:
efficiently find $\hat{f} \in F$
that fits training data

Your house: attributes

Model $\hat{f}$

Average price for
your house in CHF

# Simple linear regression in 1-d

(example with one attribute)

# **Step I**: Collect data of other houses

Input attribute $x$: size in $m^2$;

Output $y$: sales price in CHF



$x_1 = 120, y_1 = 1.5\ mil$

$x_2 = 200\ \ y_2 = 3\ mil$

$x_3 = 130\ \ y_3 = 1\ mil$
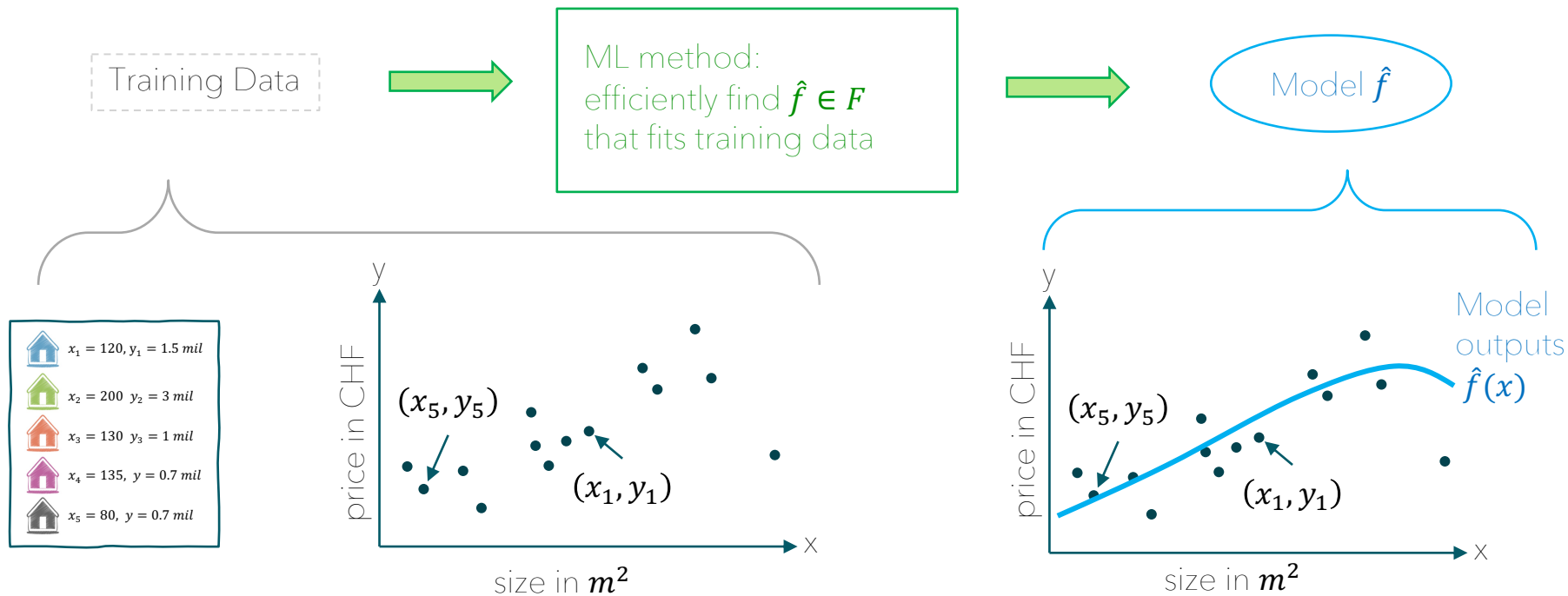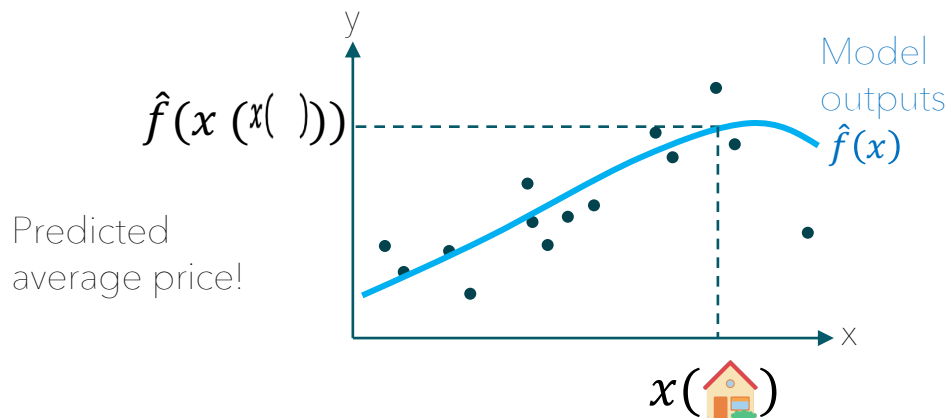
$x_4 = 135,\ y_4 = 0.7\ mil$
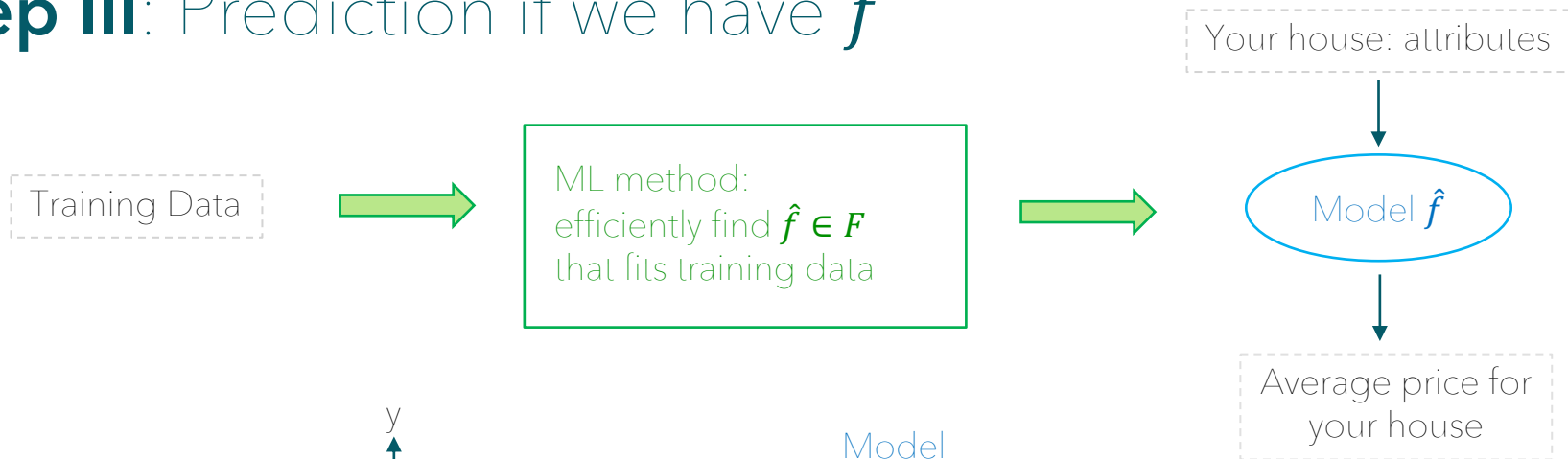
$x_5 = 80,\ y_5 = 0.7\ mil$

*Disclaimer: these are artificial numbers in an imaginary city*

# From training data to model $\hat{f}$
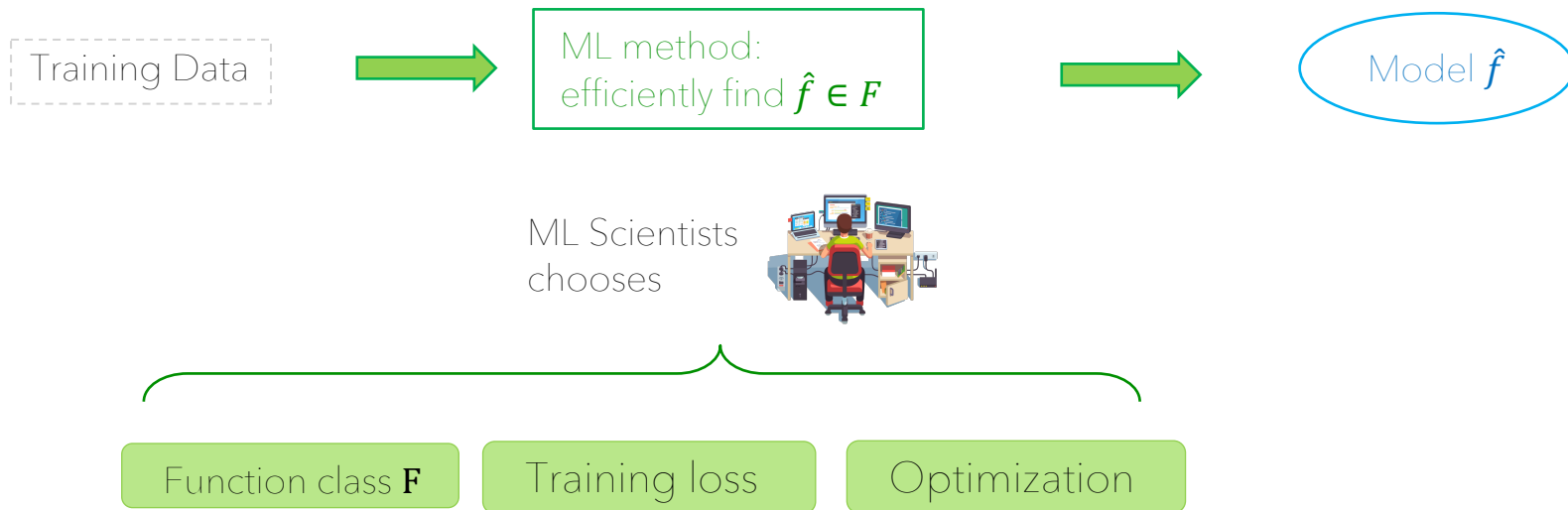
Training Data

$\longrightarrow$

ML method:
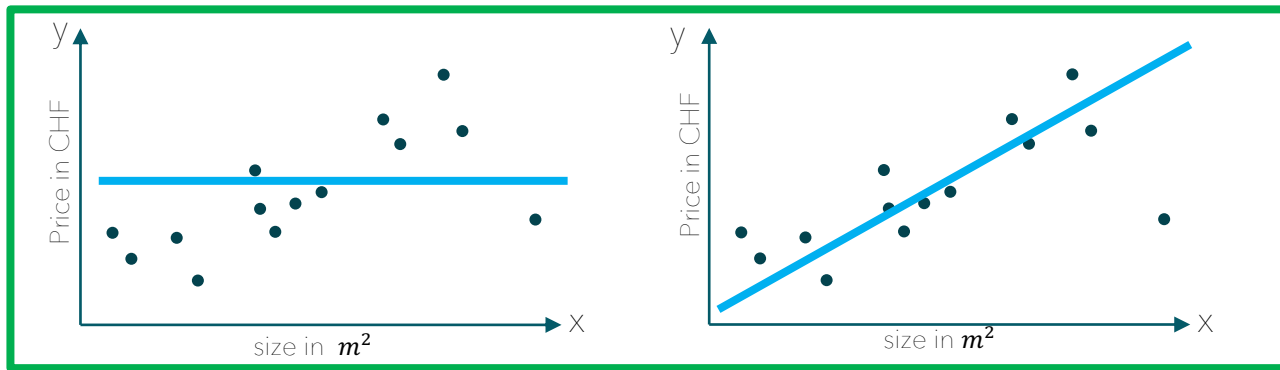efficiently find $\hat{f} \in F$
that fits training data

$\longrightarrow$

Model $\hat{f}$

$x_1 = 120, y_1 = 1.5\ mil$

$x_2 = 200\ \ y_2 = 3\ mil$

$x_3 = 130\ \ y_3 = 1\ mil$

$x_4 = 135,\ \ y = 0.7\ mil$

$x_5 = 80,\ \ y = 0.7\ mil$

y

price in CHF

$(x_5, y_5)$

$(x_1, y_1)$

size in $m^2$

X

y

price in CHF

$(x_5, y_5)$

$(x_1, y_1)$

Model
outputs
$\hat{f}(x)$

size in $m^2$

X

# **Step III**: Prediction if we have $\hat{f}$

Training Data

$\longrightarrow$

ML method:
efficiently find $\hat{f} \in F$
that fits training data

$\longrightarrow$

Your house: attributes

Model $\hat{f}$

Average price for
your house

$\hat{f}(x\,(x(\quad)))$

Model
outputs
$\hat{f}(x)$

y

Predicted
average price!

x

$x(\;🏠\;)$

# Step II: A method to obtain a model $\hat{f}$

Training Data $\rightarrow$ ML method: efficiently find $\hat{f} \in F$ $\rightarrow$ Model $\hat{f}$

ML Scientists chooses



| Function class $\mathbf{F}$ | Training loss | Optimization |

# Different function classes $F$

Question: With what kind of functions can I fit 1-dimensional data

constant
and linear

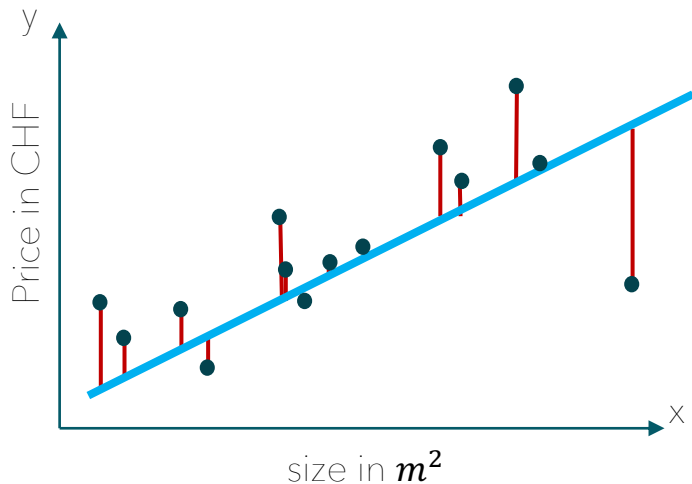more generally
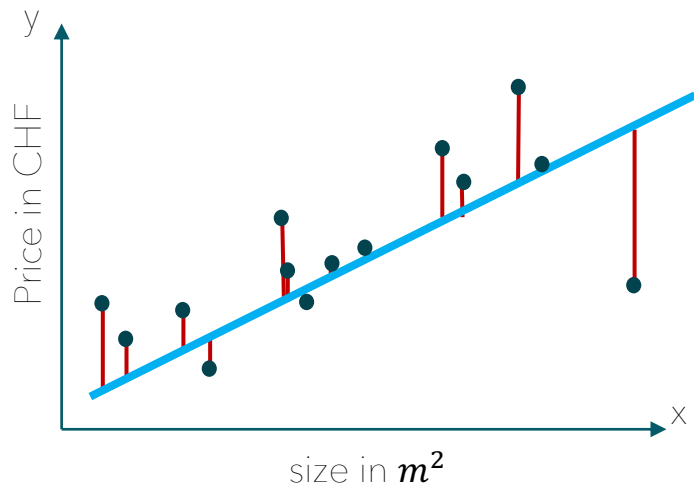polynomial, nonlinear
(next time)

# Different training losses

Question: What is a good fit of the training data?

If $f(x)$ is "close" to $y$ for most points, that is, it has low training loss $L(f) = \frac{1}{n}\sum_{i=1}^{n} \ell(f(x_i), y_i)$

Let's think about suitable loss functions, representing some notion of distance



Which loss $\ell(f(x), y)$ would you choose?

Discuss with neighbor for two minutes and answer on eduApp.

(a) $|f(x) - y|$   (b) $f(x) - y$

(c) $(f(x) - y)^2$  (d) $max(y - f(x), 0)$

# Different training losses

Question: What is a good fit of the training data?

If $f(x)$ is "close" to $y$ for most points, that is, it has low training loss $L(f) = \frac{1}{n}\sum_{i=1}^{n} \ell(f(x_i), y_i)$

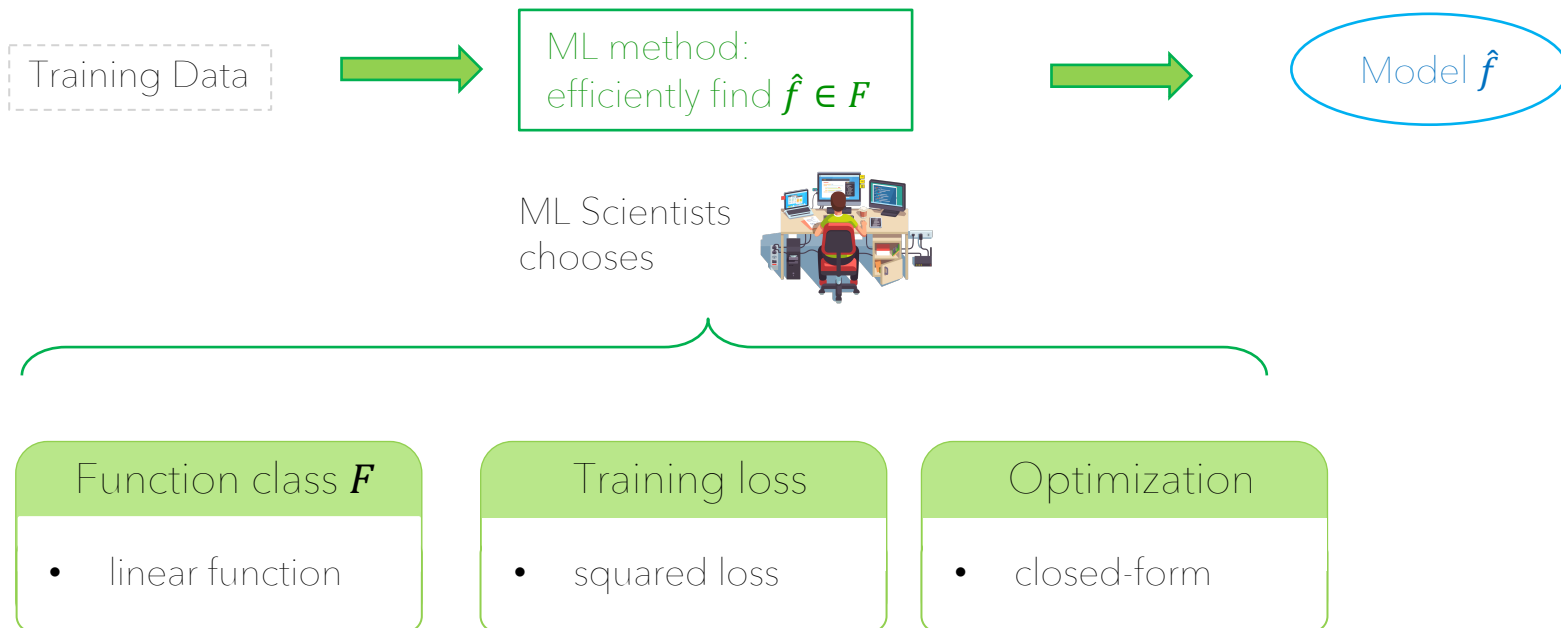Let's think about suitable loss functions, representing some notion of distance



$f(x) = -10^5$ for all $x$ would have low loss!

not differentiable

(a) $|f(x) - y|$    (b) $f(x) - y$

(c) $(f(x) - y)^2$  (d) $max(y - f(x), 0)$

overestimation is rewarded, not differentiable

18

# **Step II**: The method choice for this lecture

Training Data →

ML method:
efficiently find $\hat{f} \in F$

→ Model $\hat{f}$

ML Scientists chooses



| Function class $F$ | Training loss | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |

# Function class: Linear functions

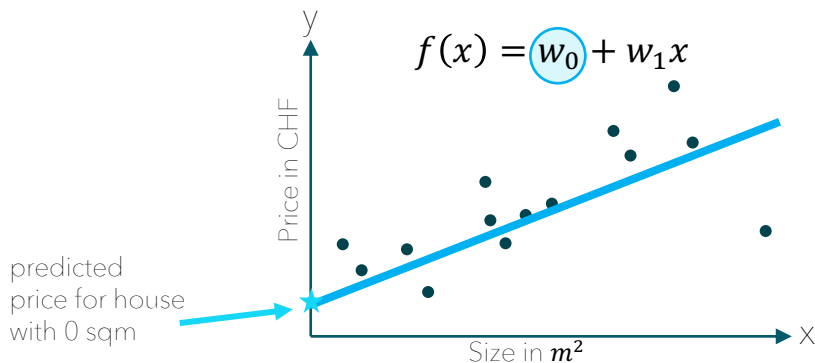| Function class $F$ | Loss function | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |

Class/set of all linear functions $F_{lin} = \{ f : f(x) = w_0 + w_1 x \text{ for } w_0, w_1 \in \mathbb{R} \}$

We say, all functions in $F_{lin}$ are *parameterized* by two scalars $w_0$, $w_1$



$f(x) = \boxed{w_0} + w_1 x$

predicted price for house with 0 sqm

y

Price in CHF

Size in $m^2$

X



$f(x) = w_0 + \boxed{w_1} x$

predicted change in price

change in $m^2$

y

Price in CHF

Size in $m^2$

X

# Training loss using the squared loss

| Function class $F$ | Training loss | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |



y — Price in CHF

x — size in $m^2$

What is a good fit of the training data?

We want to find the function $\hat{f}$ that minimizes the training loss (squared loss) on average over the training set:

$$\hat{f} = \operatorname*{argmin}_{f \in F_{lin}} L(f) = \operatorname*{argmin}_{f \in F_{lin}} \frac{1}{n} \sum_{i=1}^{n} \left(y_i - f(x_i)\right)^2$$

$\hat{f}$ we call the output of our ML model – linear regression

# Linear model $\hat{f}$ that minimizes training loss

- **Goal**: Find model $\hat{f}$ that has the smallest training loss

$$\hat{f} = \operatorname*{argmin}_{f \in F_{lin}} L(f) = \operatorname*{argmin}_{f \in F_{lin}} \frac{1}{n} \sum_{i=1}^{n} \left(y_i - f(x_i)\right)^2$$

- **Recall**: all linear functions in $F_{lin} = \{f : f(x) = w_0 + w_1 x \text{ for } w_0, w_1 \in \mathbb{R}\}$

  *are parameterized by two scalars $w_0, w_1$*

- Since $\hat{f}$ is a function in $F_{lin}$, it has to be of the form $\hat{f}(x) = \widehat{w}_0 + \widehat{w}_1 x$ for two scalars $w_0, w_1$,

  $\rightarrow$ searching for a minimum in $F_{lin}$ is the same as searching for scalars $w_0, w_1$ that minimize
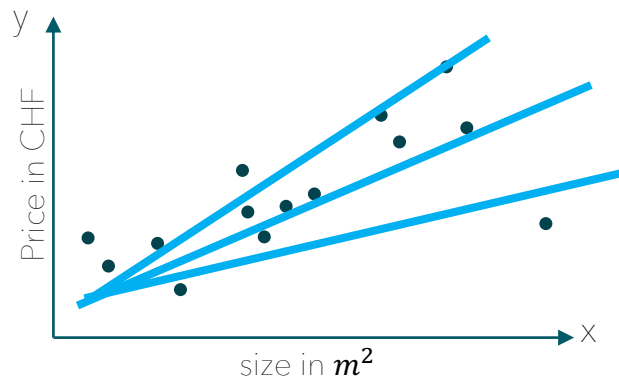
$$\widehat{w} := (\widehat{w}_0, \widehat{w}_1) = \operatorname*{argmin}_{w_0, w_1 \in \mathbb{R}} L(w_0, w_1) = \operatorname*{argmin}_{w_0, w_1 \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i)^2$$

# Minimization of the training loss: How to find $\widehat{w}$

| Function class $F$ | Loss function | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |

Training loss for linear functions $L(w_0, w_1) = \frac{1}{n}\sum_{i=1}^{n}(y_i - w_0 - w_1 x_i)^2$

We would like to find training loss minimizer $\widehat{w} = (\widehat{w}_0, \widehat{w}_1) = \underset{w_0, w_1 \in \mathbb{R}}{\operatorname{argmin}} L(w_0, w_1)$



- Some candidate linear functions (left), but $F_{lin}$ includes infinitely many candidates!

- We can compute minimizer $w_0, w_1$ analytically!

# Simplifying the problem: Fixing $w_0 = 0$

| Function class $F$ | Loss function | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |

The training loss for linear functions $L(w_0, w_1) = \frac{1}{n}\sum_{i=1}^{n}(y_i - w_0 - w_1 x_i)^2$ is a function of $w_0, w_1$

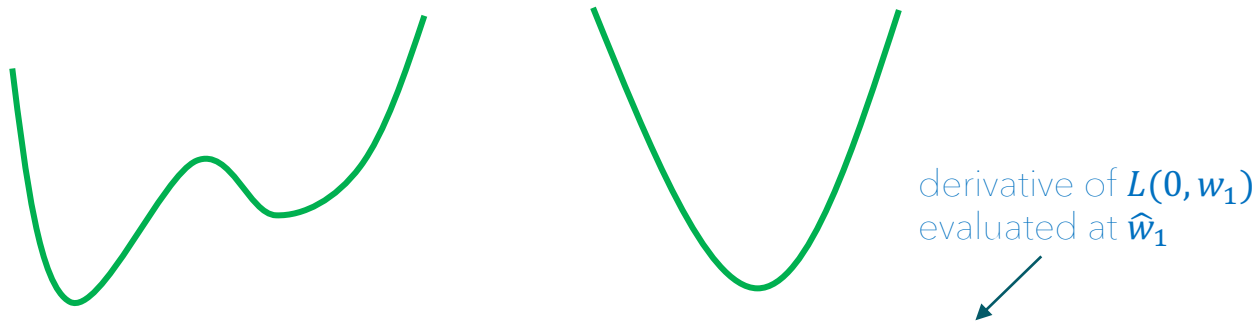For simplification, we first fix $w_0 = 0$, and only optimize $L(0, w_1)$ over $w_1 \in \mathbb{R}$

# Analysis recap: Stationary points of 1-d functions

| Function class $F$ | Loss function | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |

For a general 1-d function $g(x)$, where's the minimum $\hat{x} = \text{argmin}_{x \in R}\, g(x)$

Here are two example functions. Find all stationary points, local and global minima.



derivative of $L(0, w_1)$ evaluated at $\hat{w}_1$

→ since $L(0, w_1)$ is a one-dimensional quadratic: there's one minimum where $\boxed{L'(0, \hat{w}_1) = 0}$

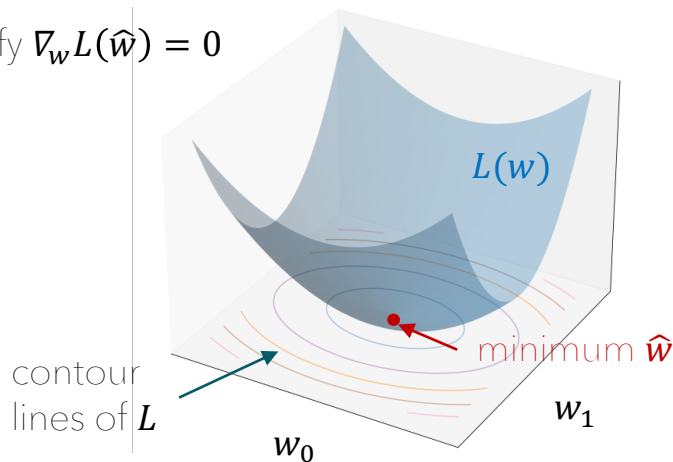# Solving the original problem: Finding the minimum $\widehat{w}$

| Function class $\mathcal{F}$ | Loss function | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |

More generally: let $w_0$ again be variable, i.e. want to find $\widehat{w} = (\widehat{w}_0, \widehat{w}_1) = \underset{w_0, w_1 \in \mathbb{R}}{\operatorname{argmin}} L(w_0, w_1)$

- By Theorem 2.3. math recap, a global minimum $\widehat{w}$ must satisfy $\nabla_w L(\widehat{w}) = 0$

- Recall the training loss $L(w_0, w_1) = \frac{1}{n}\sum_{i=1}^{n}(y_i - w_0 - w_1 x_i)^2$



$L(w)$

minimum $\widehat{w}$

contour lines of $L$

$w_0$

$w_1$

In the next five minutes: Turn to your neighbor and derive precise conditions on $\widehat{w}$ as a function of the sample points $\{(x_i, y_i)\}_{i=1}^{n}$ (calculate the gradient first)
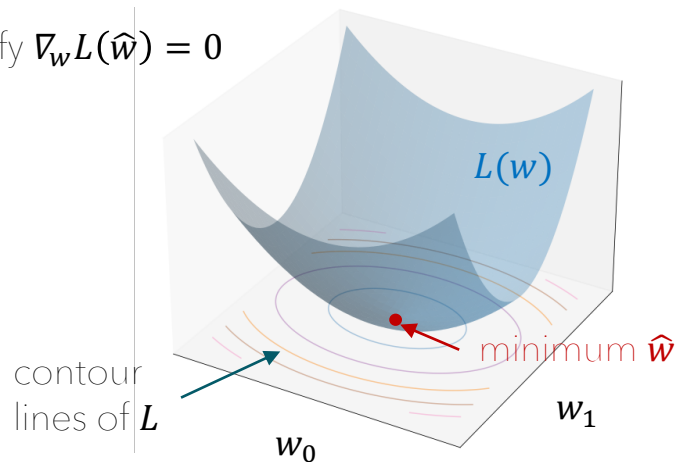
# Solving the original problem: Finding the minimum $\widehat{w}$

| Function class $F$ | Loss function | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |

More generally: let $w_0$ again be variable, i.e. want to find $\widehat{w} = (\widehat{w}_0, \widehat{w}_1) = \underset{w_0, w_1 \in \mathbb{R}}{\operatorname{argmin}} L(w_0, w_1)$

- By Theorem 2.3. math recap, a global minimum $\widehat{w}$ must satisfy $\nabla_w L(\widehat{w}) = 0$

- Recall the training loss $L(w_0, w_1) = \frac{1}{n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i)^2$

- And hence the minimum $\widehat{w} = (\widehat{w}_0, \widehat{w}_1)$ satisfies

$$\nabla_w L(\widehat{w}_0, \widehat{w}_1) = \begin{pmatrix} -\frac{2}{n} \sum_{i=1}^{n} (y_i - \widehat{w}_0 - \widehat{w}_1 x_i) \\ -\frac{2}{n} \sum_{i=1}^{n} (y_i - \widehat{w}_0 - \widehat{w}_1 x_i) \, x_i \end{pmatrix} = 0$$

$L(w)$

contour lines of $L$

minimum $\widehat{w}$

$w_1$

$w_0$

# How many minima $\widehat{w}$? – Linear Algebra refresher

A minimum $\widehat{w}$ must satisfy $\begin{pmatrix} -\frac{2}{n}\sum_{i=1}^{n}(y_i - \widehat{w}_0 - \widehat{w}_1 x_i) \\ -\frac{2}{n}\sum_{i=1}^{n}(y_i - \widehat{w}_0 - \widehat{w}_1 x_i)\, x_i \end{pmatrix} = 0.$ First think about how many

solutions exist to this equation. How many (global) minima $\widehat{w}$ exist?

(A) One          (B) Multiple          (C) None          (D) Depends on $\{(x_i, y_i)\}_{i=1}^{n}$

*Solution: (D) Regarding number of solutions:*

- *brute force: deriving closed-form solutions (see homework)*

- *or: matrix vector notation (we see later today)*

- *or: system of two linear equations for two parameters → when does it have a unique solution?*

*Further, all solutions are global minima, we'll discuss next week in more detail, why …*

# What you can do now

- high-level: teach a machine how to use training data to output a prediction

  rule/model (that can be used for prediction of a new point)

  by deriving a closed-form solution

- know what a training loss is and minimize the training loss (with square loss)

  over 1-d linear functions by minimizing over the scalars

  that parameterize the function

# Side comment: other losses

| Function class $F$ | Loss function | Optimization |
|---|---|---|
| • linear function | • different choices | • closed-form |

Squared loss

- weighs over- and underestimation the same

- Cost grows quadratically (large errors hugely penalized)

Instead might want:

- ignore outliers (ones with very large penalty) → Huber loss

- weigh over- and underestimation differently → asymmetric losses
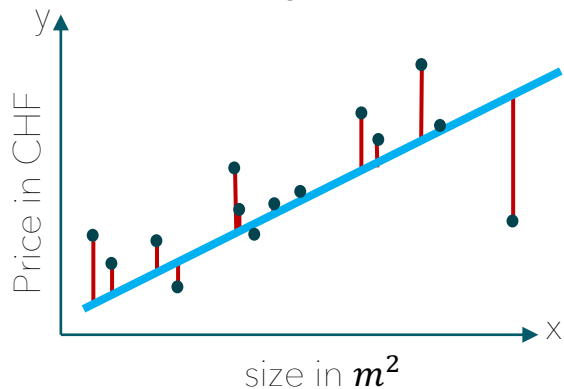


31

# Minimizer of Huber loss

| Function class $\boldsymbol{F}$ | Loss function | Optimization |
|---|---|---|
| • linear function | • Huber loss | • closed-form |

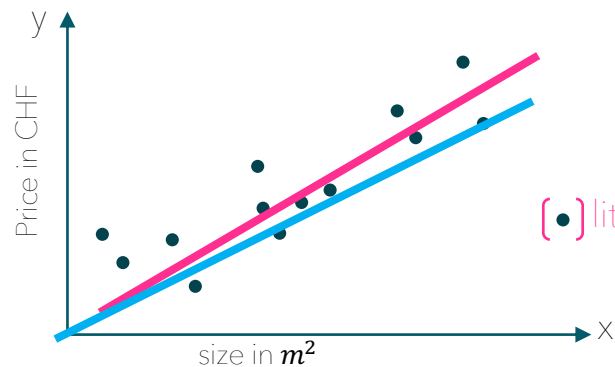Minimizer of training loss with square loss $\ell_{square}$



far points

have very little weight

Minimizer of training loss Huber loss $\ell_{Huber}$



[ • ] little weight

**Perhaps better choice if there are outliers in your training data not representing current market**

# Minimizer of asymmetric loss
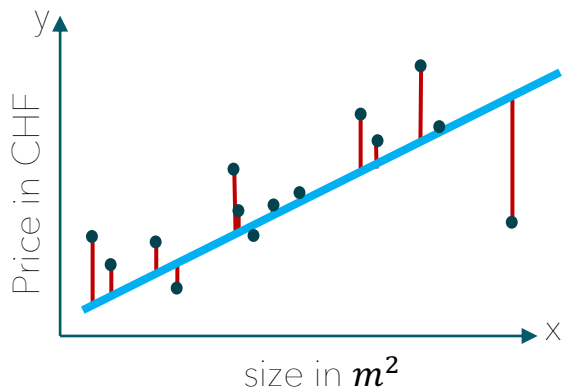
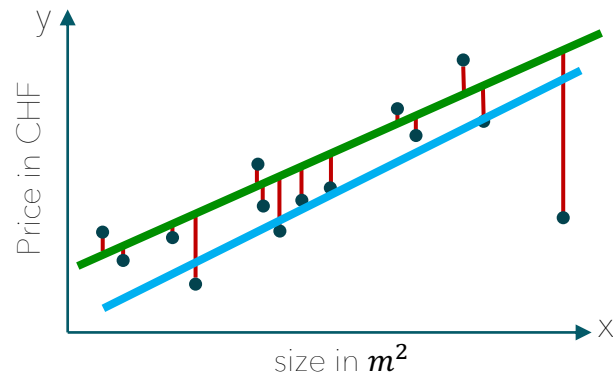| Function class $F$ | Loss function | Optimization |
|---|---|---|
| • linear function | • asymmetric regression loss | • closed-form |

Minimum of square loss $\ell_{square}$



More weight on overestimation

Minimum of loss that places less weight on overestimation $\ell_{over}$
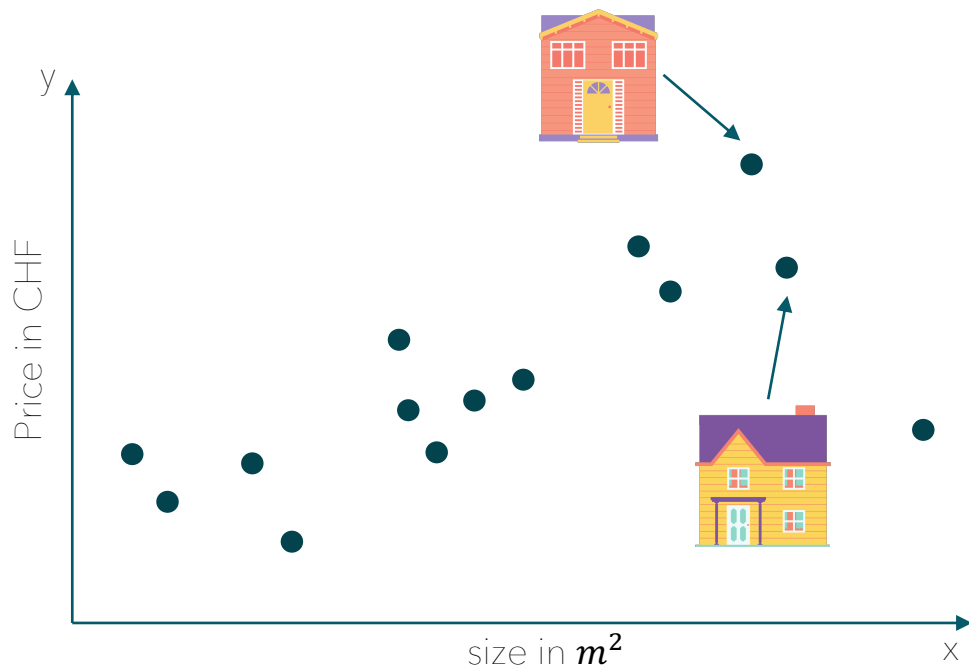


**Perhaps better choice if you care more about maximizing profit than how fast you can sell**

# Multiple regression

(Linear least squares)

# More inputs available



These two houses have similar size, but very different price?

→ this 1d model completely ignores other attributes of the house!

For example they might differ in

- number of bathrooms,
- distance to train station,
- construction year etc.

# **Step I**: Collect more data of other houses

$x_{[1]}$ = size (in $m^2$),
$x_{[2]}$ = # of bathrooms
$x_{[3]}$ = distance to train (in km),
$x_{[4]}$ = years since construction

Input: attribute vector $x = (x_{[1]}, \dots, x_{[d]})$,

Output: $y$: sales price in CHF

**extract**

$x_1 = (120,2,0.5,1), y_1 = 1.5\ mio$

$x_2 = (200,3,0.5,3)\ \ y_2 = 3\ mio$

$x_3 = (130,2,5,1)\ \ y_3 = 1\ mio$

$x_4 = (135,1,5,1)\ \ y_4 = 0.7\ mio$
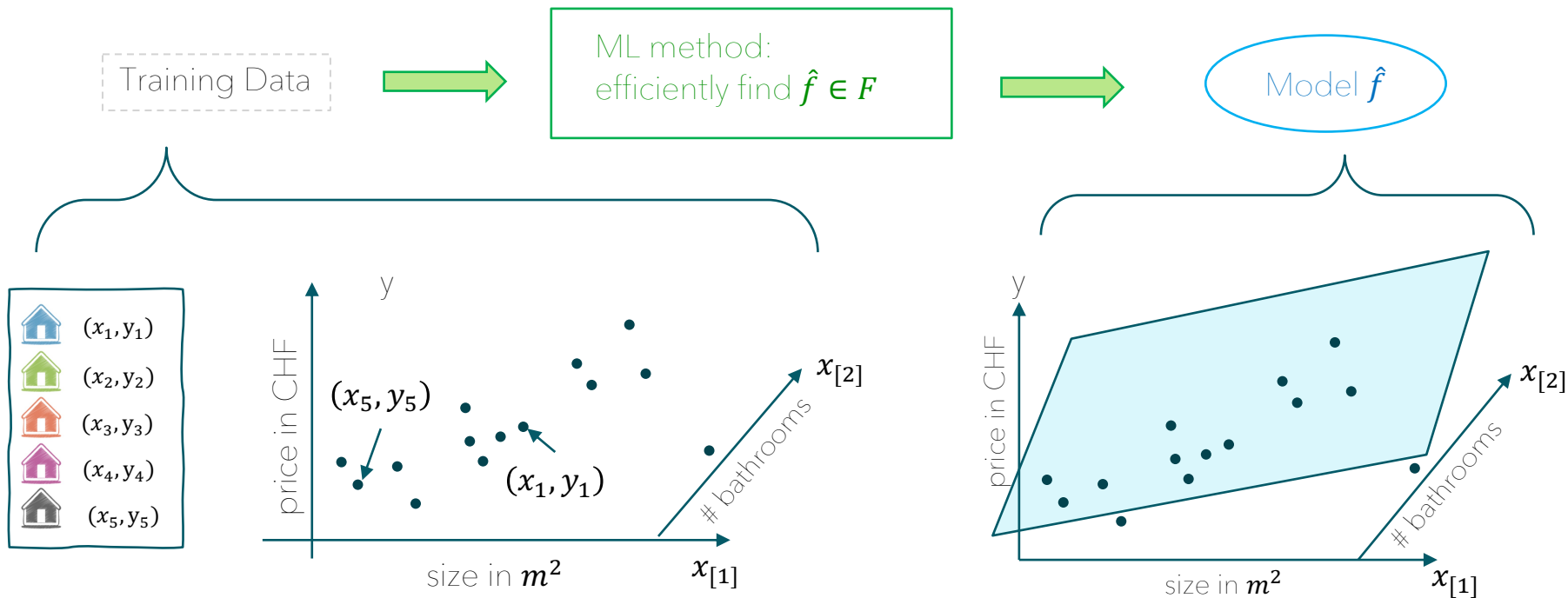
$x_5 = (80,4,1,0.5)\ \ y_5 = 0.7\ mio$

notation:
$x_i$ - input attributes of sample $i$
$x_{[i]}$ - $i$-th attribute

*Disclaimer: these are artificial numbers in an imaginary city*

# Multiple regression

(All visualizations are for two attributes, i.e. $d = 2$, the formulas are more general $d > 1$)

# Function class: Linear with $d$ linear features

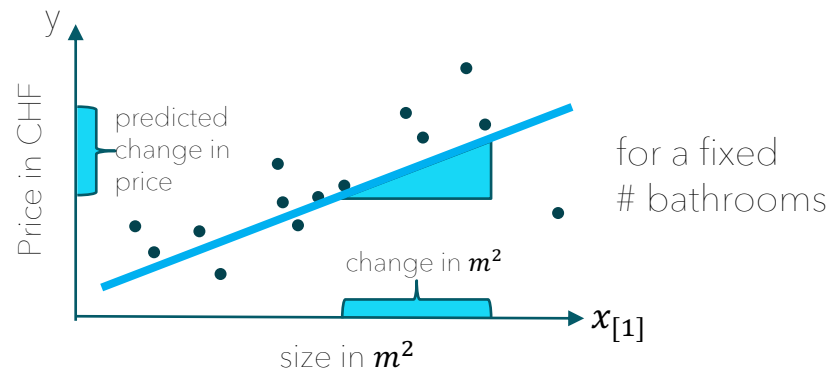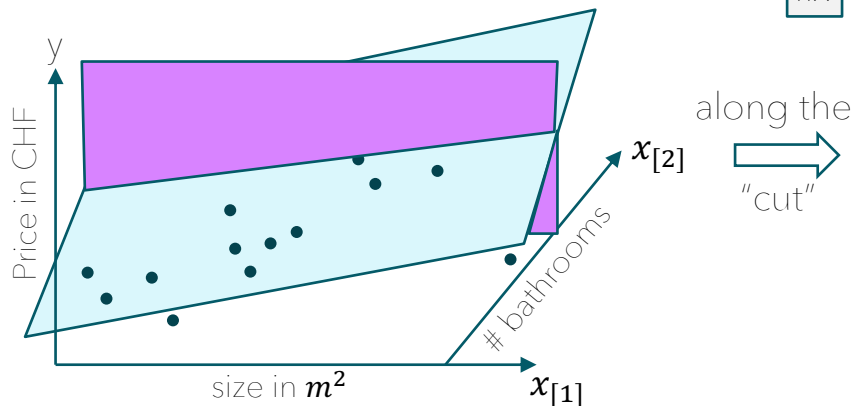| Function class $F$ | Loss function | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |

Class/set of all linear functions $F_{lin} = \{ f : f(x) = w_0 + \sum_{j=1}^{d} w_j x_{[j]} = w_0 + w^\top x \text{ for } w = (w_1, \ldots, w_d) \in \mathbb{R}^d \}$

Visualization for d = 2: $f(x) = w_0 + w_1 x_{[1]} + w_2 x_{[2]}$
fix



y

Price in CHF

size in $m^2$

# bathrooms

$x_{[2]}$

$x_{[1]}$

along the
"cut"

y

Price in CHF

predicted
change in
price

change in $m^2$

size in $m^2$

$x_{[1]}$

for a fixed
# bathrooms

# Training loss using the squared loss

| Function class $F$ | Loss function | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |



Y

W$_1$

W$_2$

Visualization for $d = 2$

Analogous to 1-d: learned model $\hat{f}$ minimizes training loss

$$\hat{f} = \operatorname*{argmin}_{f \in F_{lin}} L(f) = \operatorname*{argmin}_{f \in F_{lin}} \frac{1}{n} \sum_{i=1}^{n} \left( y_i - f(x_i) \right)^2$$

equivalent to minimizing over vector $w$

$$\hat{w} = \operatorname*{argmin}_{w_0 \in \mathbb{R},\, w \in \mathbb{R}^d} L(w_0, w) = \operatorname*{argmin}_{w_0 \in \mathbb{R},\, w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^{n} (y_i - w_0 - w^{\mathsf{T}} x_i)^2$$
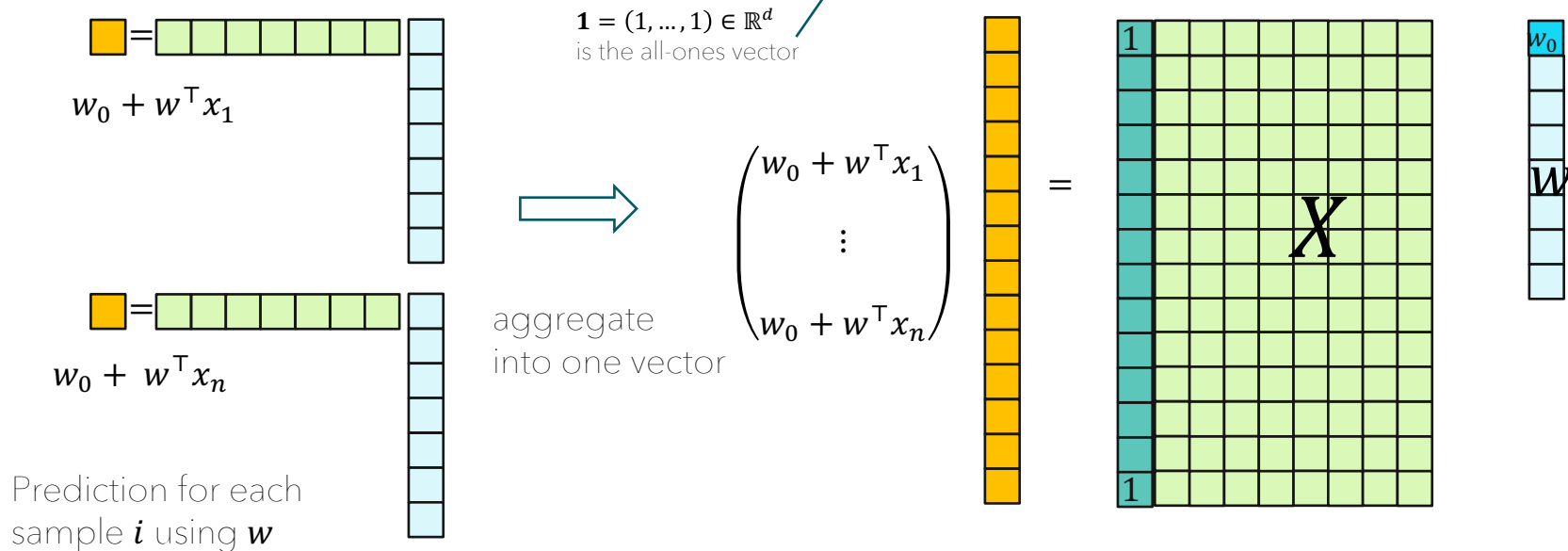
# Training loss in matrix vector notation

We now see how we can rewrite the training loss $L$ in matrix vector notation

$$L(w_0, w) = \frac{1}{n} \sum_{i=1}^{n} (y_i - w_0 - w^\top x_i)^2 = \frac{1}{n} \left\| y - \mathbf{1} w_0 - Xw \right\|^2$$



$\mathbf{1} = (1, \ldots, 1) \in \mathbb{R}^d$
is the all-ones vector

$w_0 + w^\top x_1$

$w_0 + w^\top x_n$

Prediction for each sample $i$ using $w$

aggregate into one vector

$$\begin{pmatrix} w_0 + w^\top x_1 \\ \vdots \\ w_0 + w^\top x_n \end{pmatrix} = \begin{matrix} 1 \\ \\ \\ 1 \end{matrix} \; X \; \begin{matrix} w_0 \\ \\ w \end{matrix}$$

# Minimizing the training loss

- For simplicity let's set $w_0 = 0$ and minimize over $w$

- Training loss $L(0, w) = \frac{1}{n} \left\| y - Xw \right\|^2 = \frac{1}{n} \left\| y \right\|^2 - \frac{2}{n} y^\top Xw + \frac{1}{n} w^\top X^\top Xw$

  with gradient $\nabla_w L(0, w) = \frac{2}{n} (X^\top Xw - X^\top y)$ and Hessian $D^2 L(0, w) = \frac{2}{n} X^\top X$

- We'll now look at two ways to find the minimum
  - *stationary point condition (gradient = 0)*
  - *geometric argument (orthogonal projection)*

# Optimal solution: via stationary point condition

| Function class $F$ | Loss function | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |

- Again: Minimum $\widehat{w} = \mathrm{argmin}_w\, L(0, w)$ must be a stationary point, i.e. satisfying $\nabla_w L(0, \widehat{w}) = 0$

- All stationary points of this quadratic loss $\frac{1}{n}\left\|y - Xw\right\|^2$ are minima, because the

  Hessian $\frac{2}{n} X^\top X$ is positive semi-definite (psd) *(see math recap Sec 1.6., 2.5 and details next week)*

- $\nabla_w L(0, w) = \frac{2}{n}\left(X^\top X w\ -\ X^\top y\right)$ together with stationary point condition $\nabla_w L(0, \widehat{w}) = 0$ yields

$$\rightarrow X^\top y = X^\top X\,\widehat{w} \rightarrow \widehat{w} = (X^\top X)^{-1} X^\top y$$

# Optimal solution: via geometric argument

| Function class $F$ | Loss function | Optimization |
|---|---|---|
| • linear function | • squared loss | • closed-form |

$$\widehat{w} = \text{argmin}_w \left\| \; {}^n \boxed{y} \; - \; {}^n \boxed{X}^{\; d} \; \boxed{w}^{\; d} \; \right\|$$

- set of all possible **Xw**: **span(X)** *(subspace spanned by columns of **X**, math recap)*

- let $\Pi_X$ be the orthogonal projection matrix onto **span(X)**

- then, the closest point to **y** on **span(X)** is $\Pi_X y$ (math recap Section 1.4 )



$\mathbb{R}^n$

$\|y - \Pi_X y\|$

$\hat{y} = \Pi_X y$
has smallest $\|y - \hat{y}\|$

span(X)

# Optimal solution: via geometric argument

We now derive the expression for the vector $\widehat{w}$ such that $X\widehat{w} = \Pi_X y$

We know that because $X\widehat{w}$ is an orthogonal projection

- the residual $y - X\widehat{w}$ orthogonal to all $v \in span(X) \iff (y - X\widehat{w})^\mathsf{T} Xw = 0$ for all $w \in \mathbb{R}^d$

# Optimal solution: via geometric argument

We now derive the expression for the vector $\widehat{w}$ such that $X\widehat{w} = \Pi_X y$ (see also Exerc. 1.6. in recap notes)

We know that because $X\widehat{w}$ is an orthogonal projection

- the residual $y - X\widehat{w}$ orthogonal to all $v \in span(X) \iff (y - X\widehat{w})^\top Xw = 0$ for all $w \in \mathbb{R}^d$

- hence we require $X^\top(y - X\widehat{w}) = 0 \iff X^\top y = X^\top X \widehat{w}$ (called normal equations!)

$\to$ this yields the unique solution $\widehat{w} = (X^\top X)^{-1}X^\top y$ if $X^\top X$ is invertible

Question for you until next week (Answers discussed next week):

How many minima $\widehat{w}$ does $\frac{1}{n}||y - Xw||^2$ have? Argue using the matrix $X^\top X$

(A) One (B) Multiple (C) None (D) Depends on $\{(x_i, y_i)\}_{i=1}^n$

# What you can do now

On a high level:

- for inputs with multiple attributes

- teach a machine how to use training data to output a prediction rule/model

  (that can be used for prediction of a new point)

- by deriving a closed-form solution for the best linear fit

  (in the square loss sense) on the training data (training loss minimzer)

# References / acknowledgements

- We will provide lecture notes that explain the derivations in more detail

- Other reference: ISLR Chapter 3