

6.7900: Machine Learning

Lecture 4

Lecture start: Tues/Thurs 2:35pm

Who's speaking today? Prof. Tamara Broderick

Course website: gradml.mit.edu

Questions? Ask here or on piazza.com/mit/fall2024/67900/

Materials: Slides, video, etc linked from gradml.mit.edu after the lecture (but there is no livestream)

Last Times

- I. Empirical risk minimization (ERM)
- II. Maximum likelihood estimate (MLE)
- III. No features

Today's Plan

- I. MLE & ERM for supervised learning
- II. Linear regression: Why and how
- III. Some challenges

Why linear regression?

Why linear regression?

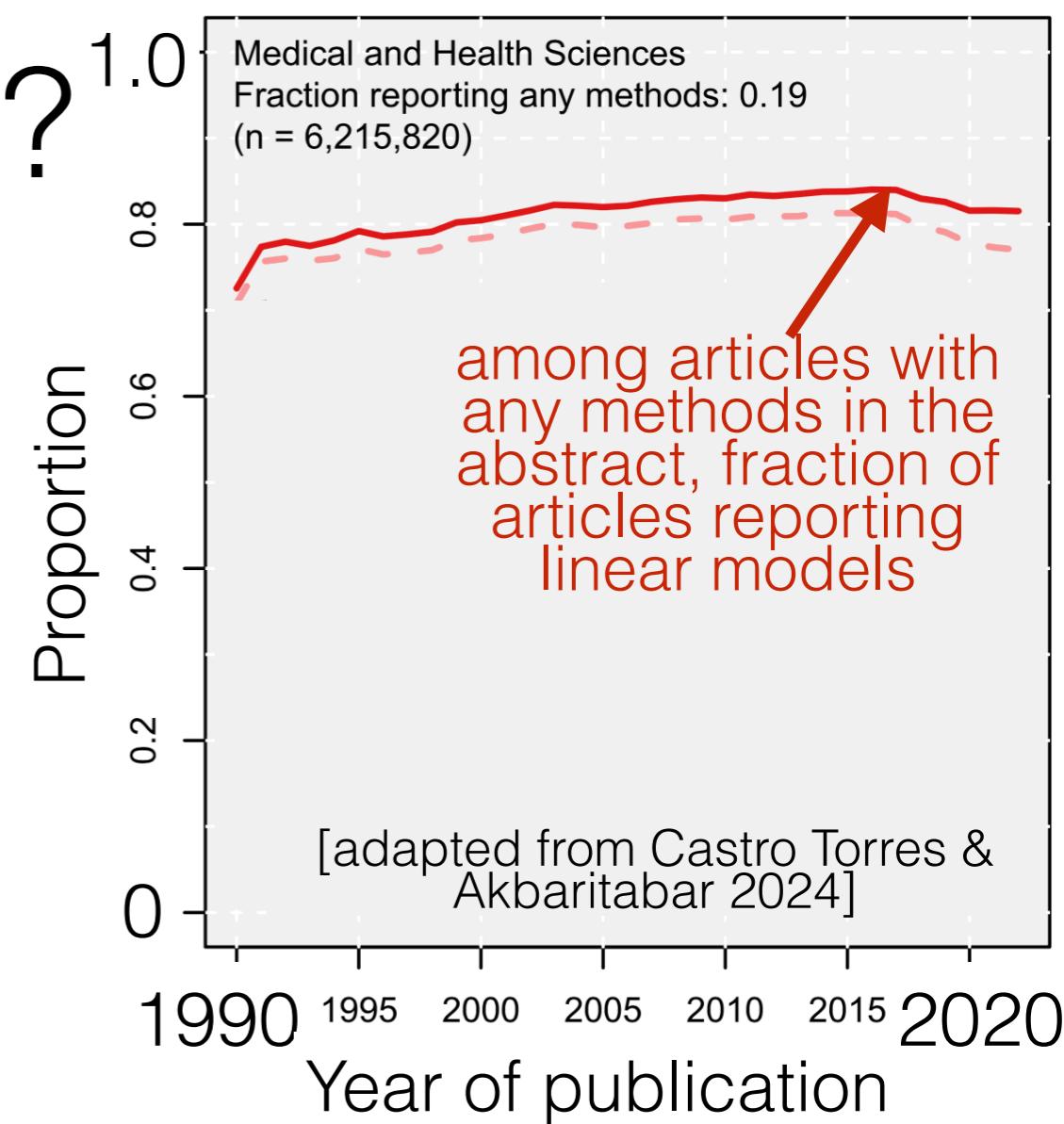
- Do modern, sophisticated robots make hammers obsolete?

Why linear regression?

- Do modern, sophisticated robots make hammers obsolete?
- Widely used

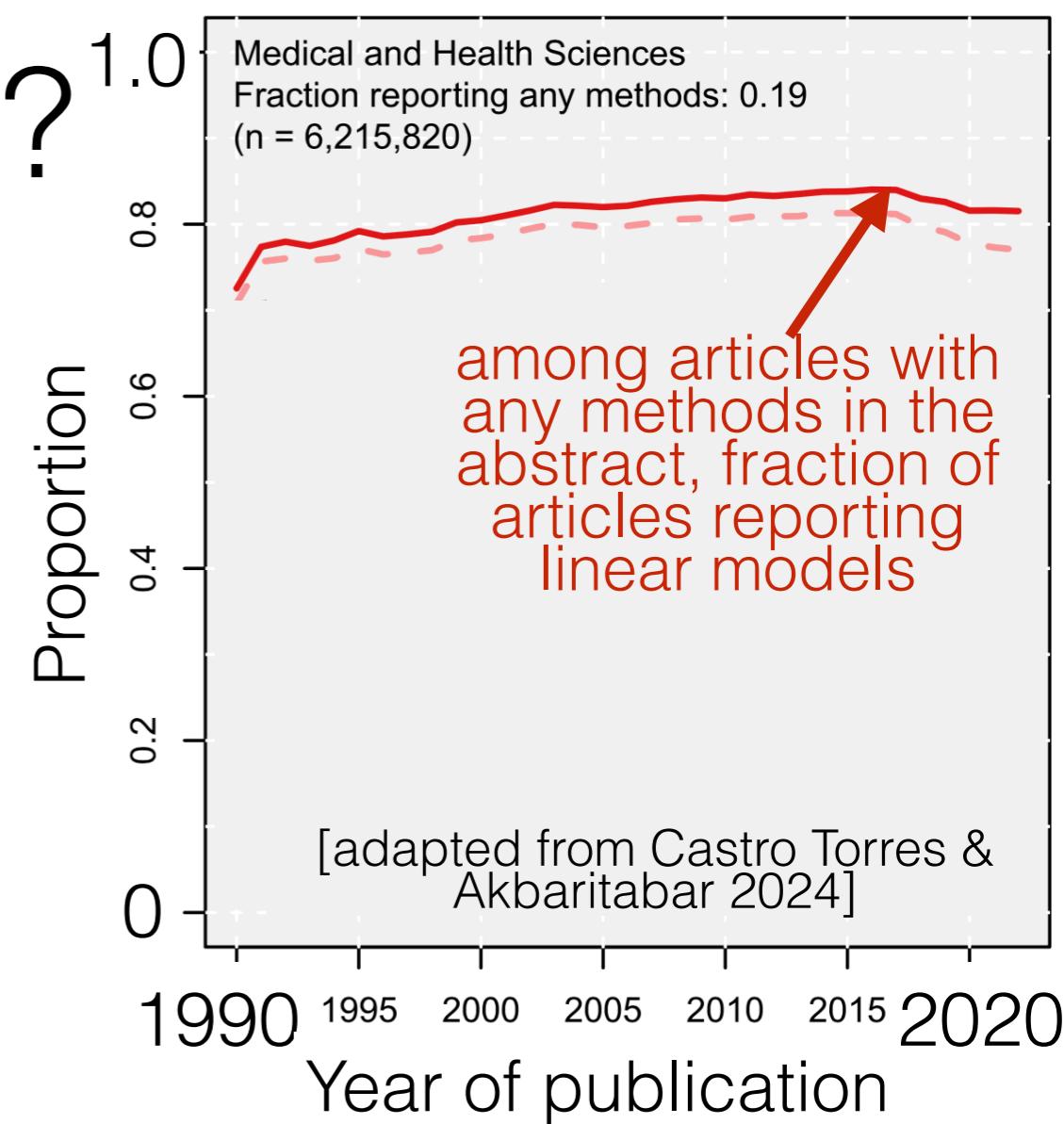
Why linear regression?

- Do modern, sophisticated robots make hammers obsolete?
- Widely used



Why linear regression?

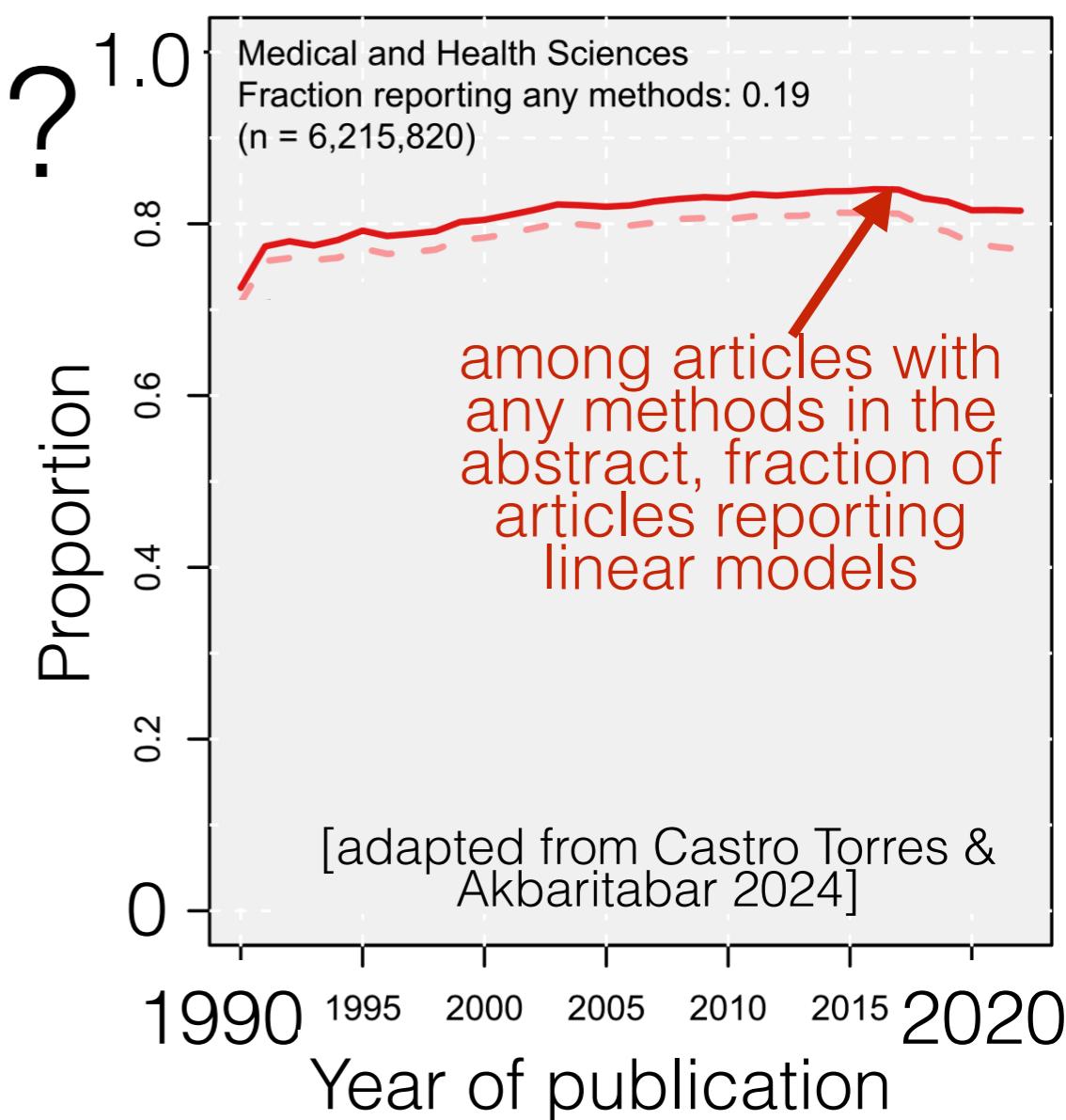
- Do modern, sophisticated robots make hammers obsolete?
- Widely used
- Speed, interpretability



Why linear regression?

- Do modern, sophisticated robots make hammers obsolete?
- Widely used
- Speed, interpretability
- Theoretical insight, including about deep learning

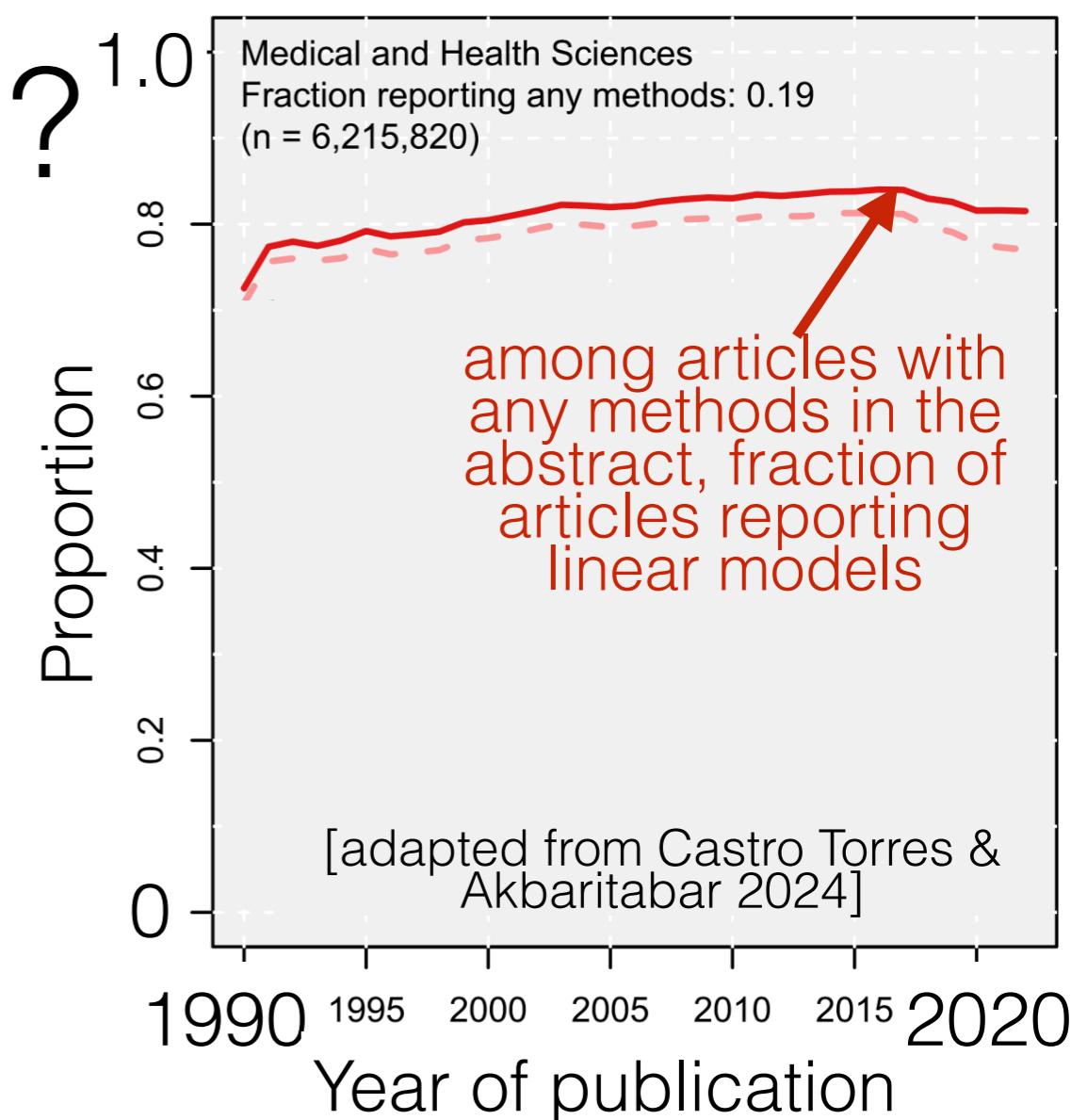
[e.g. Bartlett+
2020, Benign
overfitting in
linear
regression]



Why linear regression?

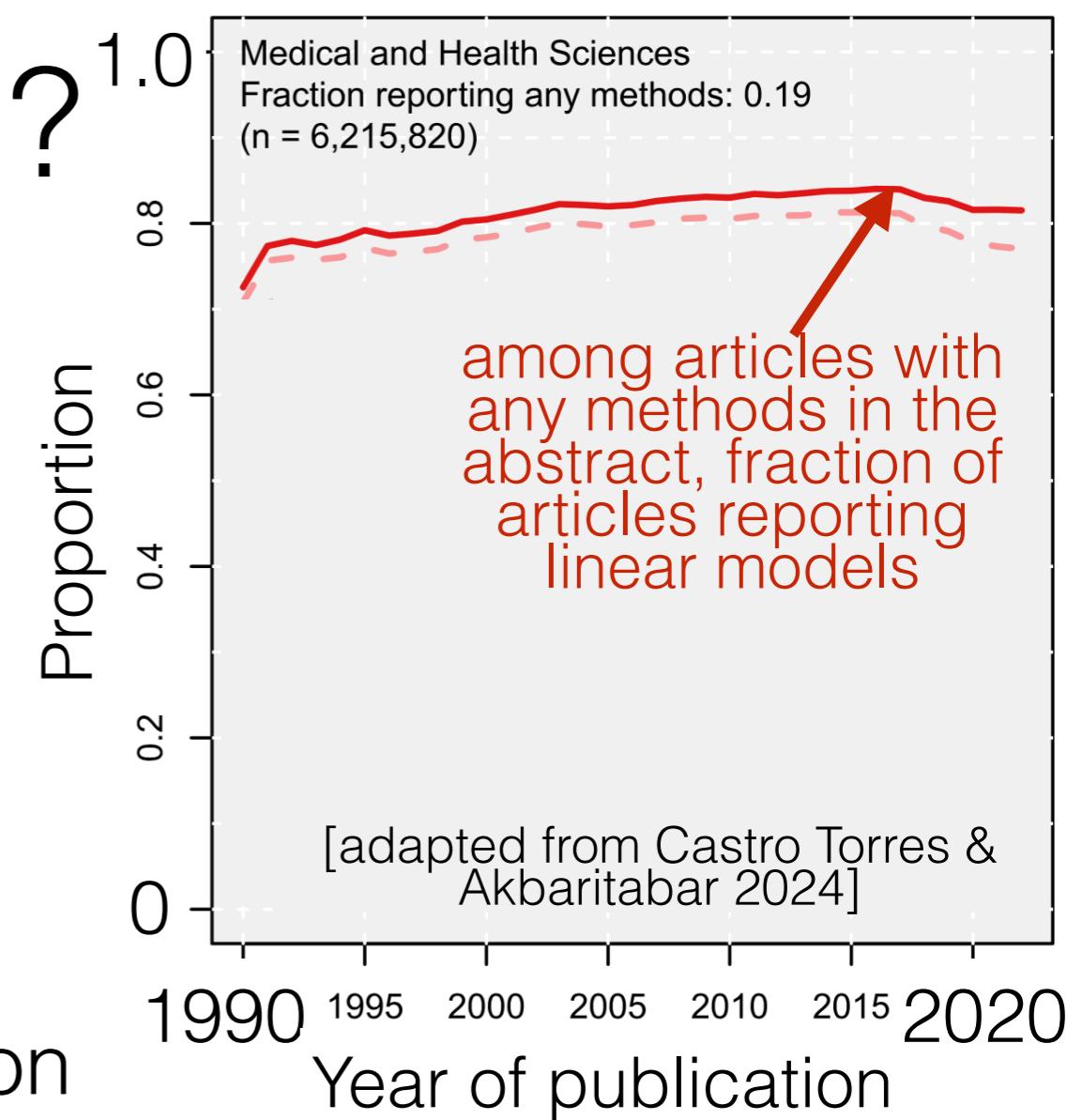
- Do modern, sophisticated robots make hammers obsolete?
- Widely used
- Speed, interpretability
- Theoretical insight, including about deep learning
- Often a cog in modern procedures

[e.g. Bartlett+
2020, Benign
overfitting in
linear
regression]



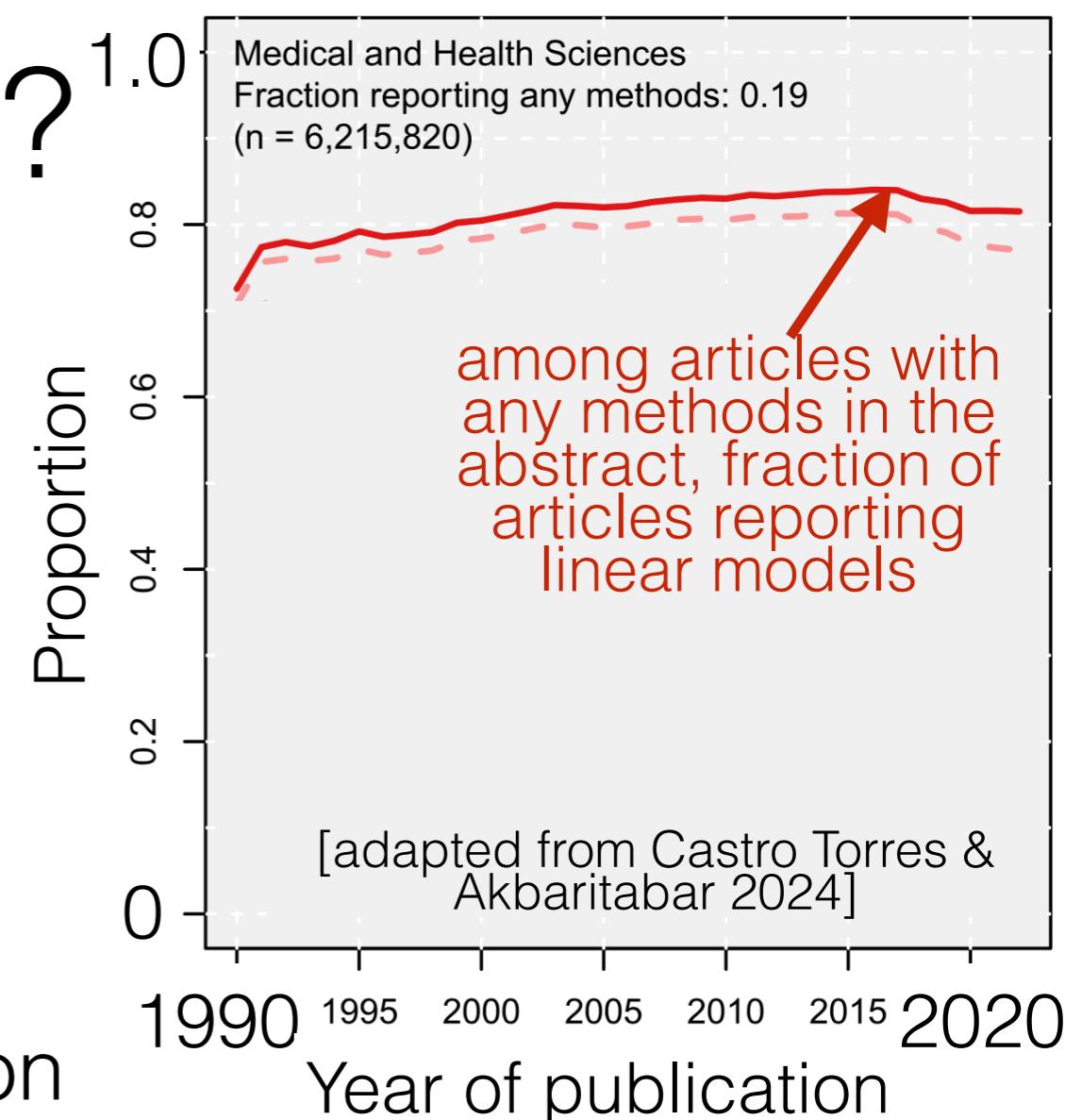
Why linear regression?

- Do modern, sophisticated robots make hammers obsolete?
- Widely used
- Speed, interpretability
- Theoretical insight, including about deep learning [e.g. Bartlett+ 2020, Benign overfitting in linear regression]
- Often a cog in modern procedures
- Importance of baselines and ablation



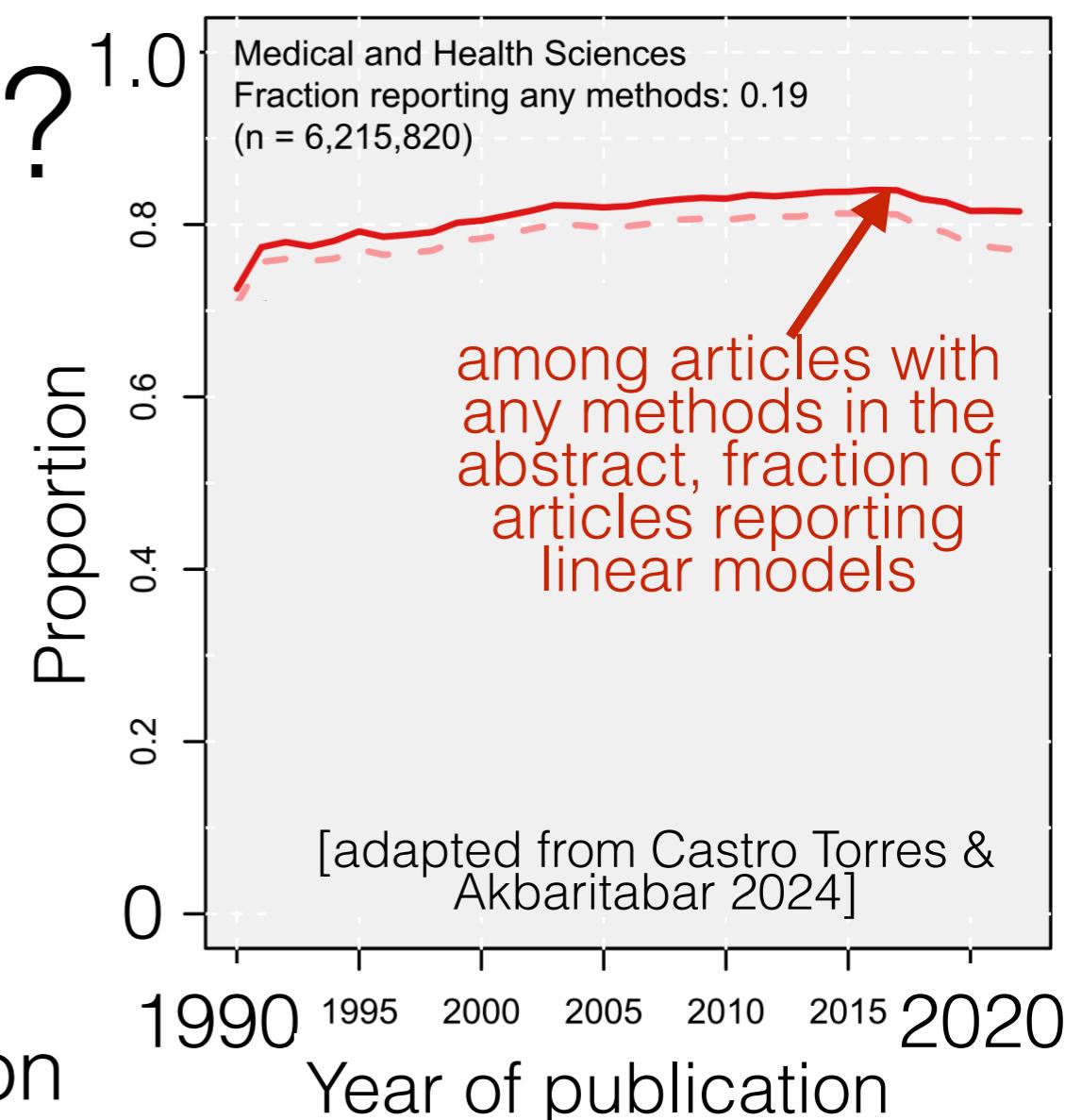
Why linear regression?

- Do modern, sophisticated robots make hammers obsolete?
- Widely used
- Speed, interpretability
- Theoretical insight, including about deep learning [e.g. Bartlett+ 2020, Benign overfitting in linear regression]
- Often a cog in modern procedures
- Importance of baselines and ablation
 - Lipton, Steinhardt 2019, Troubling Trends in Machine Learning Scholarship: Some ML papers suffer from flaws that could mislead the public and stymie future research



Why linear regression?

- Do modern, sophisticated robots make hammers obsolete?
- Widely used
- Speed, interpretability
- Theoretical insight, including about deep learning [e.g. Bartlett+ 2020, Benign overfitting in linear regression]
- Often a cog in modern procedures
- Importance of baselines and ablation
 - Lipton, Steinhardt 2019, Troubling Trends in Machine Learning Scholarship: Some ML papers suffer from flaws that could mislead the public and stymie future research
 - Salganik+ 2020: “100s of researchers attempted to predict 6 life outcomes, such as a child’s GPA and whether a family would be evicted[....]hey drew on a vast dataset that was painstakingly collected [...] over 15y. [N]o one made very accurate predictions.”

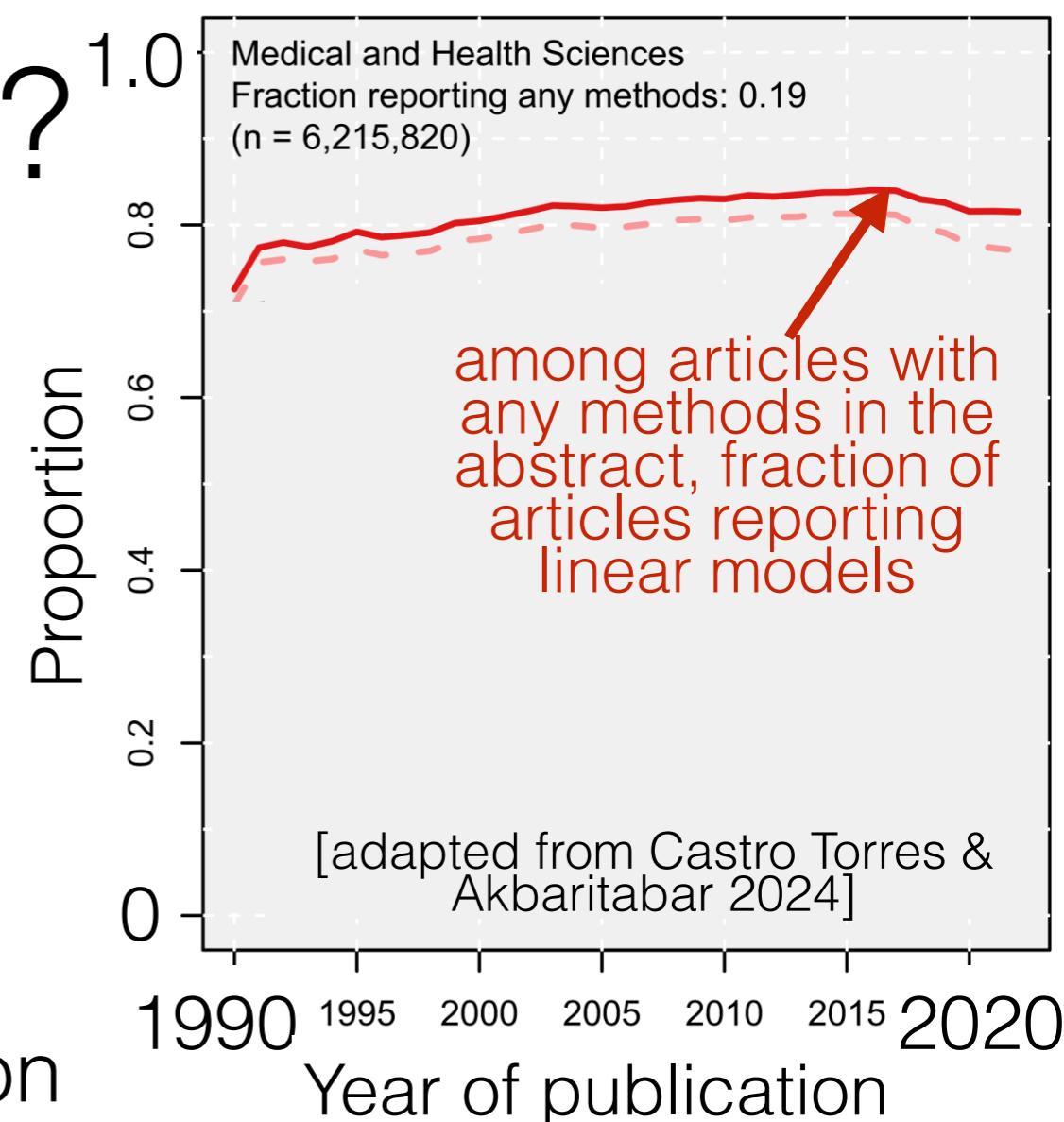


Why linear regression?

- Do modern, sophisticated robots make hammers obsolete?
- Widely used
- Speed, interpretability
- Theoretical insight, including about deep learning
- Often a cog in modern procedures

[e.g. Bartlett+ 2020, Benign overfitting in linear regression]

- Importance of baselines and ablation
 - Lipton, Steinhardt 2019, Troubling Trends in Machine Learning Scholarship: Some ML papers suffer from flaws that could mislead the public and stymie future research
 - Salganik+ 2020: “100s of researchers attempted to predict 6 life outcomes, such as a child’s GPA and whether a family would be evicted[....]hey drew on a vast dataset that was painstakingly collected [...] over 15y. [N]o one made very accurate predictions.”
 - Chen et al 2016, A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task



Motivating example

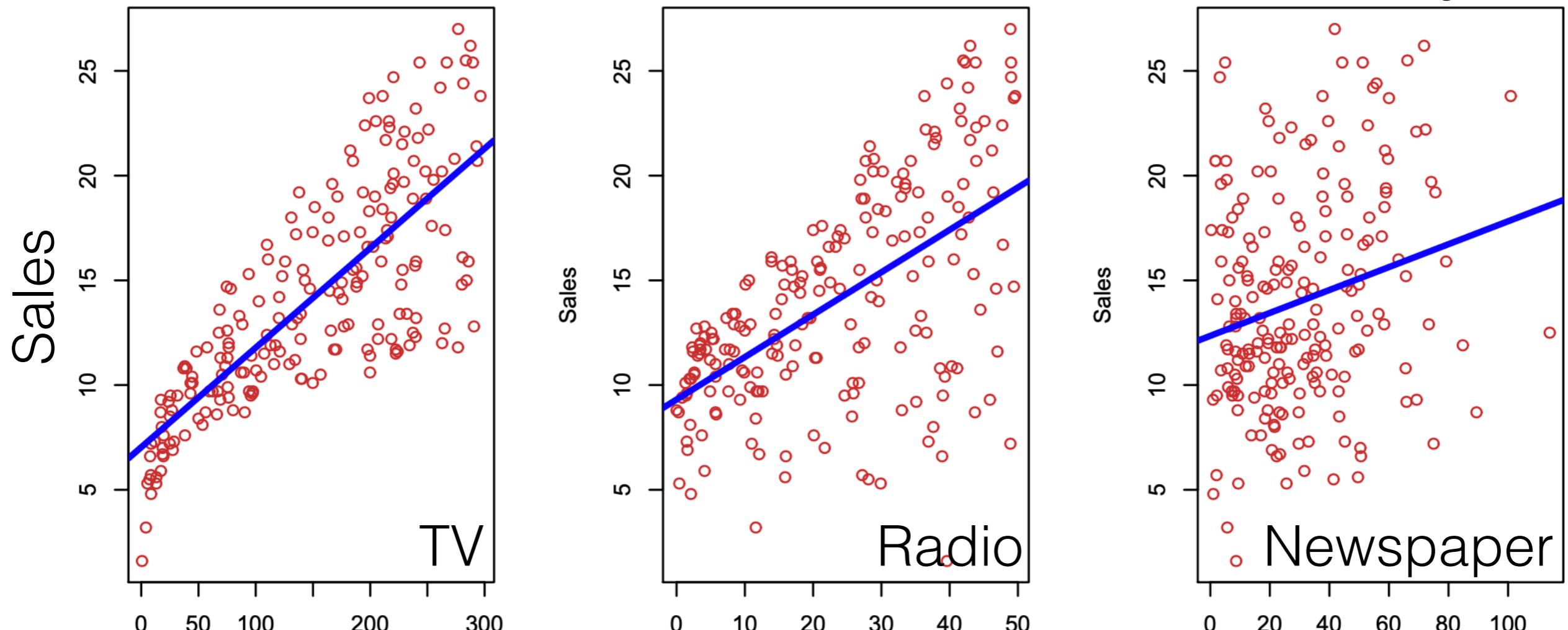
Motivating example

- Advertising data: “sales, in thousands of units, as a function of TV, radio, and newspaper budgets, in thousands of dollars, for 200 different markets”

Motivating example

- Advertising data: “sales, in thousands of units, as a function of TV, radio, and newspaper budgets, in thousands of dollars, for 200 different markets”

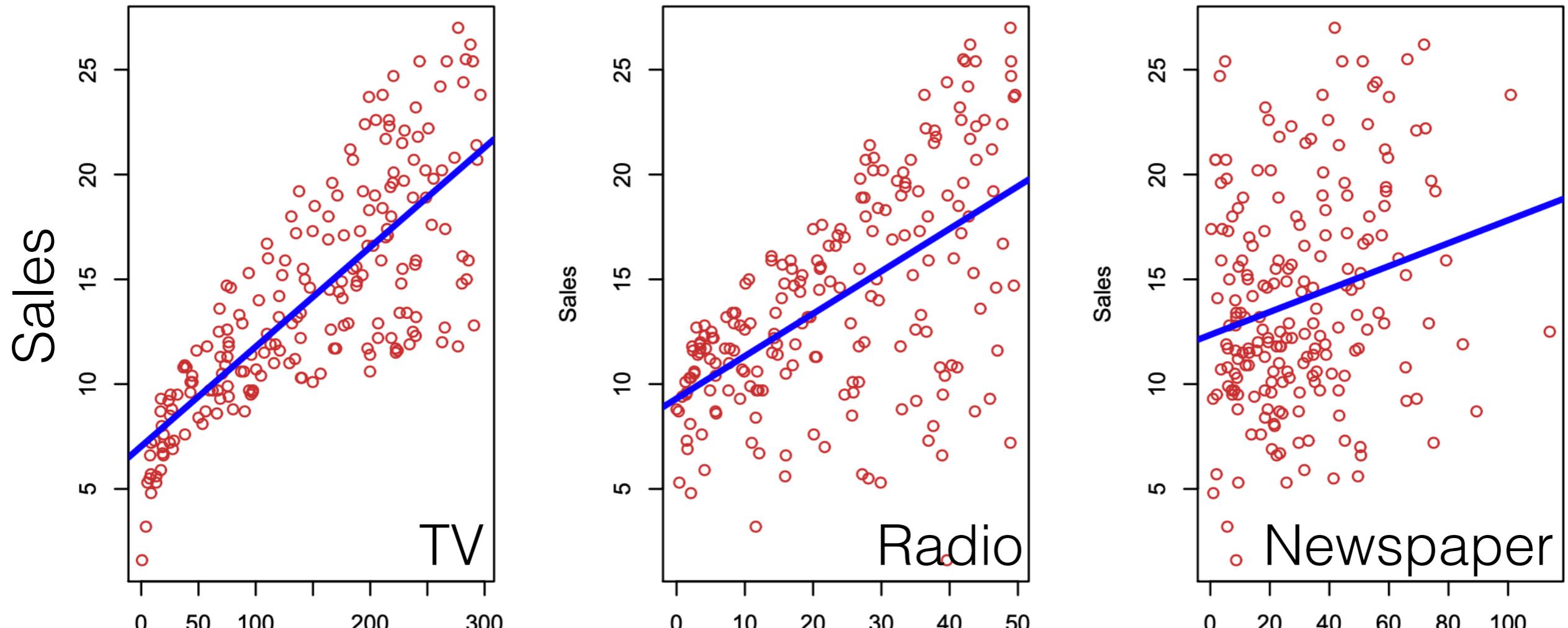
[An Intro to Statistical Learning 2023, Fig 2.1]



Motivating example

- Advertising data: “sales, in thousands of units, as a function of TV, radio, and newspaper budgets, in thousands of dollars, for 200 different markets”

[An Intro to Statistical Learning 2023, Fig 2.1]

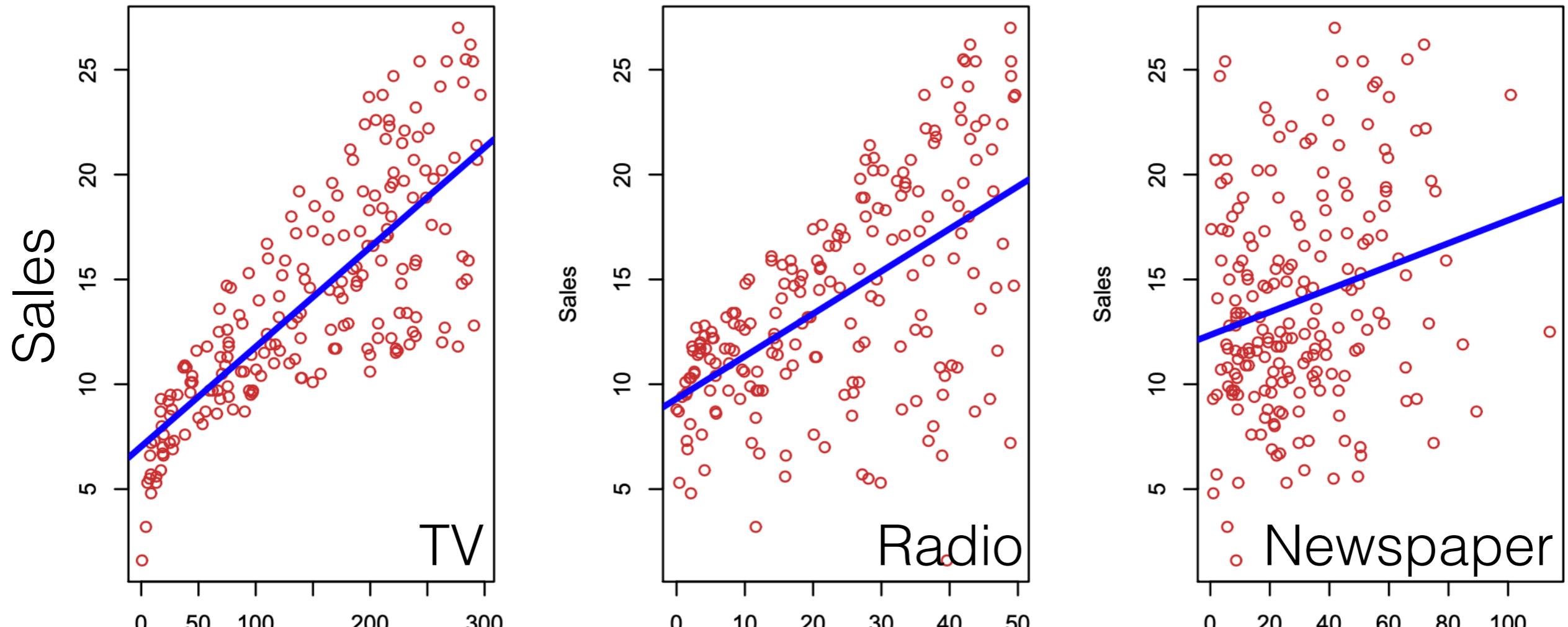


- Data point: label is sales, in thousands of units: $y^{(n)} \in \mathbb{R}_+$
- Features: $x^{(n)} = [x_1^{(n)}, x_2^{(n)}, x_3^{(n)}]^\top \in \mathbb{R}_+^3$

Motivating example

- Advertising data: “sales, in thousands of units, as a function of TV, radio, and newspaper budgets, in thousands of dollars, for 200 different markets”

[An Intro to Statistical Learning 2023, Fig 2.1]



- Data point: label is sales, in thousands of units: $y^{(n)} \in \mathbb{R}_+$
- Features: $x^{(n)} = [x_1^{(n)}, x_2^{(n)}, x_3^{(n)}]^\top \in \mathbb{R}_+^3$
- Want to predict sales under different ad budgets, allocate ad budget. Separate regressions not same as whole.

Choosing a likelihood

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}, \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$
 - Likelihood: $y^{(n)}_{1 \times 1} = \theta^\top x^{(n)}_{D \times 1} + \epsilon^{(n)}, \epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$
1x1 1xD Dx1

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$
1x1 1xD Dx1 1x1

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$

“All models
are wrong”
What about
this one?

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ homoskedastic
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ noise

“All models
are wrong”
What about
this one?

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ homoskedastic
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ noise
 - Can handle an intercept term by taking: $\forall n, x_1^{(n)} = 1$

"All models
are wrong"
What about
this one?

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ homoskedastic
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ noise
 - Can handle an intercept term by taking: $\forall n, x_1^{(n)} = 1$
 - Equivalently:

“All models
are wrong”
What about
this one?

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ **homoskedastic** “All models are wrong”
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ **noise** What about **this one?**
 - Can handle an intercept term by taking: $\forall n, x_1^{(n)} = 1$
 - Equivalently: $y^{(n)} | x^{(n)} \stackrel{indep}{\sim} \mathcal{N}(\theta^\top x^{(n)}, \sigma^2)$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ **homoskedastic** “All models are wrong”
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ **noise** What about **this one?**
 - Can handle an intercept term by taking: $\forall n, x_1^{(n)} = 1$
 - Equivalently: $y^{(n)} | x^{(n)} \stackrel{indep}{\sim} \mathcal{N}(\theta^\top x^{(n)}, \sigma^2)$
 $p(y|x, \theta, \sigma^2) = \mathcal{N}(y|\theta^\top x, \sigma^2)$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x,y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ homoskedastic
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ noise
 - Can handle an intercept term by taking: $\forall n, x_1^{(n)} = 1$
 - Equivalently: $y^{(n)} | x^{(n)} \stackrel{indep}{\sim} \mathcal{N}(\theta^\top x^{(n)}, \sigma^2)$
$$p(y|x, \theta, \sigma^2) = \mathcal{N}(y|\theta^\top x, \sigma^2)$$
 - Taking the MLE approach we've described so far:

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x, y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ homoskedastic noise
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ All models are wrong What about this one?
 - Can handle an intercept term by taking: $\forall n, x_1^{(n)} = 1$
 - Equivalently: $y^{(n)} | x^{(n)} \stackrel{\text{indep}}{\sim} \mathcal{N}(\theta^\top x^{(n)}, \sigma^2)$
$$p(y|x, \theta, \sigma^2) = \mathcal{N}(y|\theta^\top x, \sigma^2)$$
- Taking the MLE approach we've described so far:

$$\text{fit params} = \arg \max_{\text{parameters}} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \text{parameters?})$$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x, y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ **homoskedastic** **"All models are wrong"**
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ **noise** **What about this one?**
 - Can handle an intercept term by taking: $\forall n, x_1^{(n)} = 1$
 - Equivalently: $y^{(n)} | x^{(n)} \stackrel{indep}{\sim} \mathcal{N}(\theta^\top x^{(n)}, \sigma^2)$
$$p(y|x, \theta, \sigma^2) = \mathcal{N}(y|\theta^\top x, \sigma^2)$$
 - Taking the MLE approach we've described so far:

$$\hat{\theta}, \hat{\sigma}^2 = \arg \max_{\theta, \sigma^2} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2)$$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x, y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ homoskedastic noise
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ “All models are wrong” What about this one?
 - Can handle an intercept term by taking: $\forall n, x_1^{(n)} = 1$
 - Equivalently: $y^{(n)} | x^{(n)} \stackrel{\text{indep}}{\sim} \mathcal{N}(\theta^\top x^{(n)}, \sigma^2)$
$$p(y|x, \theta, \sigma^2) = \mathcal{N}(y|\theta^\top x, \sigma^2)$$

$$\hat{\theta}, \hat{\sigma}^2 = \arg \max_{\theta, \sigma^2} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) \text{ & use } p(y|x, \hat{\theta}, \hat{\sigma}^2)$$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x, y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ homoskedastic noise
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ “All models are wrong” What about this one?
 - Can handle an intercept term by taking: $\forall n, x_1^{(n)} = 1$
 - Equivalently: $y^{(n)} | x^{(n)} \stackrel{\text{indep}}{\sim} \mathcal{N}(\theta^\top x^{(n)}, \sigma^2)$
$$p(y|x, \theta, \sigma^2) = \mathcal{N}(y|\theta^\top x, \sigma^2)$$
 - Taking the MLE approach we've described so far:
$$\hat{\theta}, \hat{\sigma}^2 = \arg \max_{\theta, \sigma^2} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) \quad \& \text{use } p(y|x, \hat{\theta}, \hat{\sigma}^2)$$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x, y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ homoskedastic noise
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ “All models are wrong” What about this one?
 - Can handle an intercept term by taking: $\forall n, x_1^{(n)} = 1$
 - Equivalently: $y^{(n)} | x^{(n)} \stackrel{\text{indep}}{\sim} \mathcal{N}(\theta^\top x^{(n)}, \sigma^2)$
$$p(y|x, \theta, \sigma^2) = \mathcal{N}(y|\theta^\top x, \sigma^2)$$
 - Taking the MLE approach we've described so far:
$$\hat{\theta}, \hat{\sigma}^2 = \arg \max_{\theta, \sigma^2} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) \quad \& \text{use } p(y|x, \hat{\theta}, \hat{\sigma}^2)$$

Choosing a likelihood

Recall proposition
from Lecture 1

- **Proposition.** Consider regression with $\mathcal{X} = \mathbb{R}^D, \mathcal{Y} = \mathbb{R}$
 - Assume X, Y density $p(x, y)$ & square loss $L(a, g) = (a - g)^2$
 - (Assume all needed integrals exist.) Then this decision rule minimizes the risk of a new point: $h(x) = \mathbb{E}[Y|X = x]$
- One approach: we can make a model for $p(y|x)$ and then use the maximum likelihood estimate of the parameter
 - Due to the form of $h(x)$, it seems like we don't need to worry about a model for $p(x)$ homoskedastic
 - Likelihood: $y^{(n)} = \theta^\top x^{(n)} + \epsilon^{(n)}$, $\epsilon^{(n)} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ noise “All models are wrong” What about this one?
 - Can handle an intercept term by taking: $\forall n, x_1^{(n)} = 1$
 - Equivalently: $y^{(n)} | x^{(n)} \stackrel{indep}{\sim} \mathcal{N}(\theta^\top x^{(n)}, \sigma^2)$
$$p(y|x, \theta, \sigma^2) = \mathcal{N}(y|\theta^\top x, \sigma^2)$$
 - Taking the MLE approach we've described so far:
$$\hat{\theta}, \hat{\sigma}^2 = \arg \max_{\theta, \sigma^2} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) \text{ & use } p(y|x, \hat{\theta}, \hat{\sigma}^2)$$
 - So here predict $h(x) = \mathbb{E}[Y|X = x] = \hat{\theta}^\top x$ (no variance param!)

Maximizing the likelihood

Maximizing the likelihood

- Likelihood of the training data:

$$p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2)$$

Maximizing the likelihood

- Likelihood of the training data:

$$p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) = \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta, \sigma^2)$$

Maximizing the likelihood

- Likelihood of the training data:

$$\begin{aligned} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta, \sigma^2) \\ &= \prod_{n=1}^N \mathcal{N}(y^{(n)} | \theta^\top x^{(n)}, \sigma^2) \end{aligned}$$

Maximizing the likelihood

- Likelihood of the training data:

$$\begin{aligned} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta, \sigma^2) \\ &= \prod_{n=1}^N \mathcal{N}(y^{(n)} | \theta^\top x^{(n)}, \sigma^2) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y^{(n)} - \theta^\top x^{(n)})^2}{2\sigma^2} \right\} \end{aligned}$$

Maximizing the likelihood

- Likelihood of the training data:

$$\begin{aligned} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta, \sigma^2) \\ &= \prod_{n=1}^N \mathcal{N}(y^{(n)} | \theta^\top x^{(n)}, \sigma^2) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y^{(n)} - \theta^\top x^{(n)})^2}{2\sigma^2} \right\} \end{aligned}$$

- Equivalent to minimize the negative log likelihood (NLL)

Maximizing the likelihood

- Likelihood of the training data:

$$\begin{aligned} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta, \sigma^2) \\ &= \prod_{n=1}^N \mathcal{N}(y^{(n)} | \theta^\top x^{(n)}, \sigma^2) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y^{(n)} - \theta^\top x^{(n)})^2}{2\sigma^2} \right\} \end{aligned}$$

- Equivalent to minimize the negative log likelihood (NLL)

- $-\log p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2)$

Maximizing the likelihood

- Likelihood of the training data:

$$\begin{aligned} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta, \sigma^2) \\ &= \prod_{n=1}^N \mathcal{N}(y^{(n)} | \theta^\top x^{(n)}, \sigma^2) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y^{(n)} - \theta^\top x^{(n)})^2}{2\sigma^2} \right\} \end{aligned}$$

- Equivalent to minimize the negative log likelihood (NLL)

$$\begin{aligned} -\log p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= (N/2) \log \sigma^2 + (2\sigma^2)^{-1} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2 + \text{const} \end{aligned}$$

Maximizing the likelihood

- Likelihood of the training data:

$$\begin{aligned} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta, \sigma^2) \\ &= \prod_{n=1}^N \mathcal{N}(y^{(n)} | \theta^\top x^{(n)}, \sigma^2) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y^{(n)} - \theta^\top x^{(n)})^2}{2\sigma^2} \right\} \end{aligned}$$

- Equivalent to minimize the negative log likelihood (NLL)

$$\begin{aligned} -\log p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= (N/2) \log \sigma^2 + (2\sigma^2)^{-1} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2 + \text{const} \end{aligned}$$

- The optimal $\hat{\theta}$ can be found without finding $\hat{\sigma}^2$

Maximizing the likelihood

- Likelihood of the training data:

$$\begin{aligned} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta, \sigma^2) \\ &= \prod_{n=1}^N \mathcal{N}(y^{(n)} | \theta^\top x^{(n)}, \sigma^2) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y^{(n)} - \theta^\top x^{(n)})^2}{2\sigma^2} \right\} \end{aligned}$$

- Equivalent to minimize the negative log likelihood (NLL)

$$\begin{aligned} -\log p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= (N/2) \log \sigma^2 + (2\sigma^2)^{-1} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2 + \text{const} \end{aligned}$$

- The optimal $\hat{\theta}$ can be found without finding $\hat{\sigma}^2$

- $\hat{\theta}$ minimizes $\sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2$

Maximizing the likelihood

- Likelihood of the training data:

$$\begin{aligned} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta, \sigma^2) \\ &= \prod_{n=1}^N \mathcal{N}(y^{(n)} | \theta^\top x^{(n)}, \sigma^2) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y^{(n)} - \theta^\top x^{(n)})^2}{2\sigma^2} \right\} \end{aligned}$$

- Equivalent to minimize the negative log likelihood (NLL)

$$\begin{aligned} -\log p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= (N/2) \log \sigma^2 + (2\sigma^2)^{-1} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2 + \text{const} \end{aligned}$$

- The optimal $\hat{\theta}$ can be found without finding $\hat{\sigma}^2$

- $\hat{\theta}$ minimizes $\sum_{n=1}^N \underbrace{(y^{(n)} - \theta^\top x^{(n)})^2}_{\text{residual (signed)}}$

Maximizing the likelihood

- Likelihood of the training data:

$$\begin{aligned} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta, \sigma^2) \\ &= \prod_{n=1}^N \mathcal{N}(y^{(n)} | \theta^\top x^{(n)}, \sigma^2) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y^{(n)} - \theta^\top x^{(n)})^2}{2\sigma^2} \right\} \end{aligned}$$

- Equivalent to minimize the negative log likelihood (NLL)

$$\begin{aligned} -\log p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= (N/2) \log \sigma^2 + (2\sigma^2)^{-1} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2 + \text{const} \end{aligned}$$

- The optimal $\hat{\theta}$ can be found without finding $\hat{\sigma}^2$

- $\hat{\theta}$ minimizes $\sum_{n=1}^N \underbrace{(y^{(n)} - \theta^\top x^{(n)})^2}_{\text{residual (signed)}}$

- The **residual** is a more general concept: $y^{(n)} - h(x^{(n)})$

Maximizing the likelihood

- Likelihood of the training data:

$$\begin{aligned} p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta, \sigma^2) \\ &= \prod_{n=1}^N \mathcal{N}(y^{(n)} | \theta^\top x^{(n)}, \sigma^2) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y^{(n)} - \theta^\top x^{(n)})^2}{2\sigma^2} \right\} \end{aligned}$$

- Equivalent to minimize the negative log likelihood (NLL)

$$\begin{aligned} -\log p(\{y^{(n)}\}_1^N | \{x^{(n)}\}_1^N, \theta, \sigma^2) &= (N/2) \log \sigma^2 + (2\sigma^2)^{-1} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2 + \text{const} \end{aligned}$$

- The optimal $\hat{\theta}$ can be found without finding $\hat{\sigma}^2$

- $\hat{\theta}$ minimizes $\sum_{n=1}^N \underbrace{(y^{(n)} - \theta^\top x^{(n)})^2}_{\text{residual (signed)}}$

residual (signed)

 Residual sum of squares (RSS)

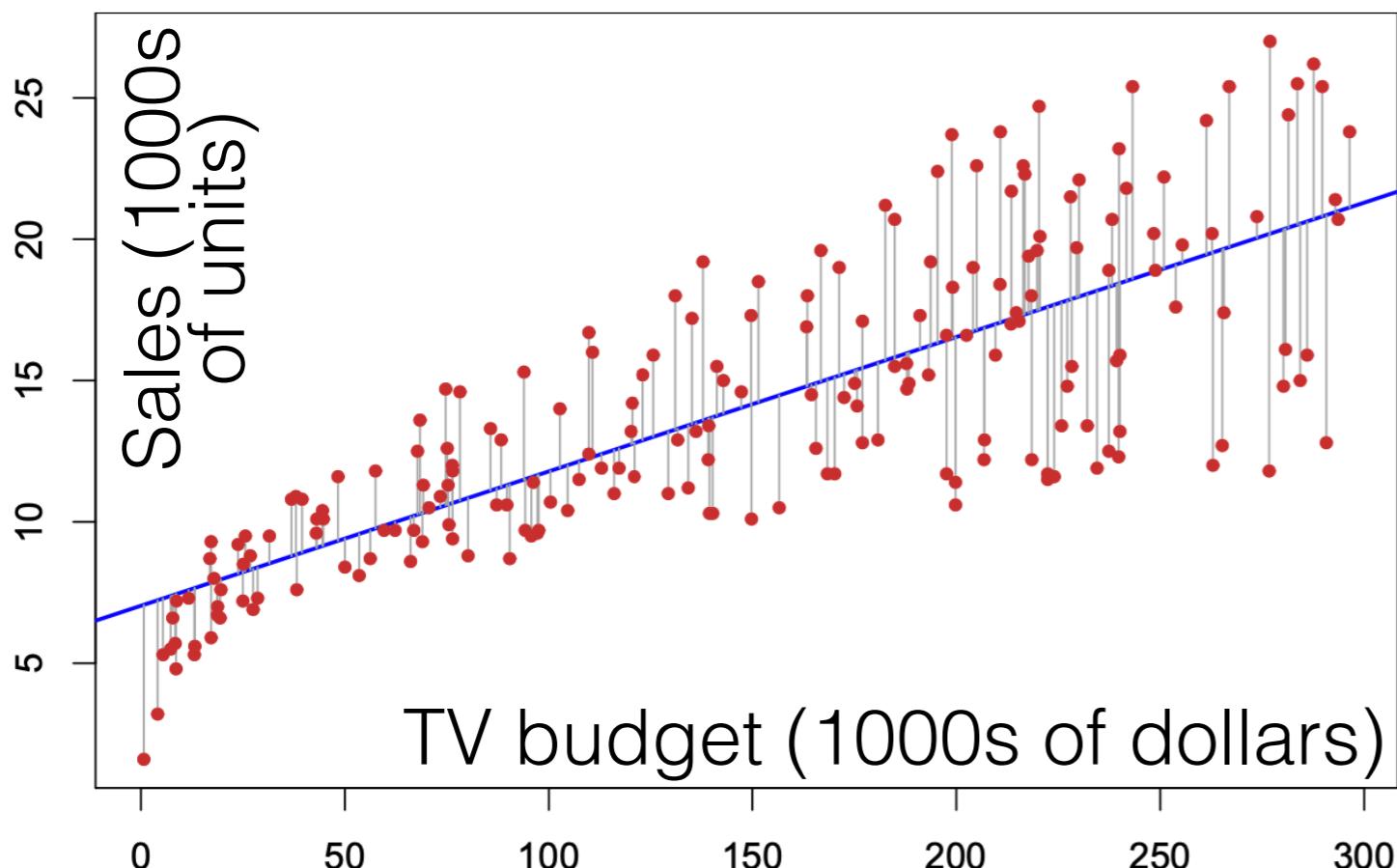
- The **residual** is a more general concept: $y^{(n)} - h(x^{(n)})$

Some visualizations

Some visualizations

- Linear predictor and residuals for one feature (with an intercept)

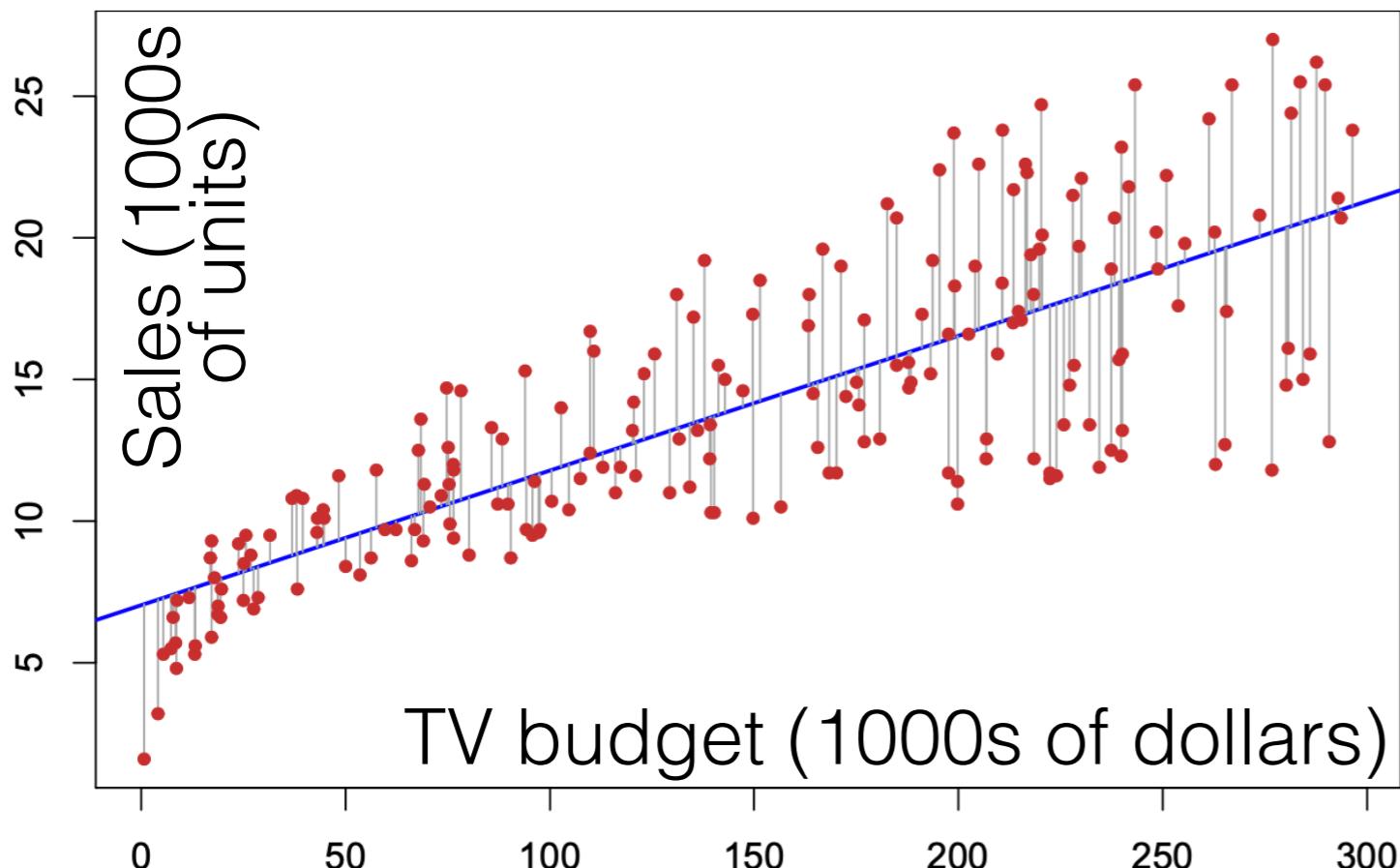
Some visualizations



- Linear predictor and residuals for one feature (with an intercept)

[An Intro to Statistical Learning 2023,
Fig 3.1]

Some visualizations

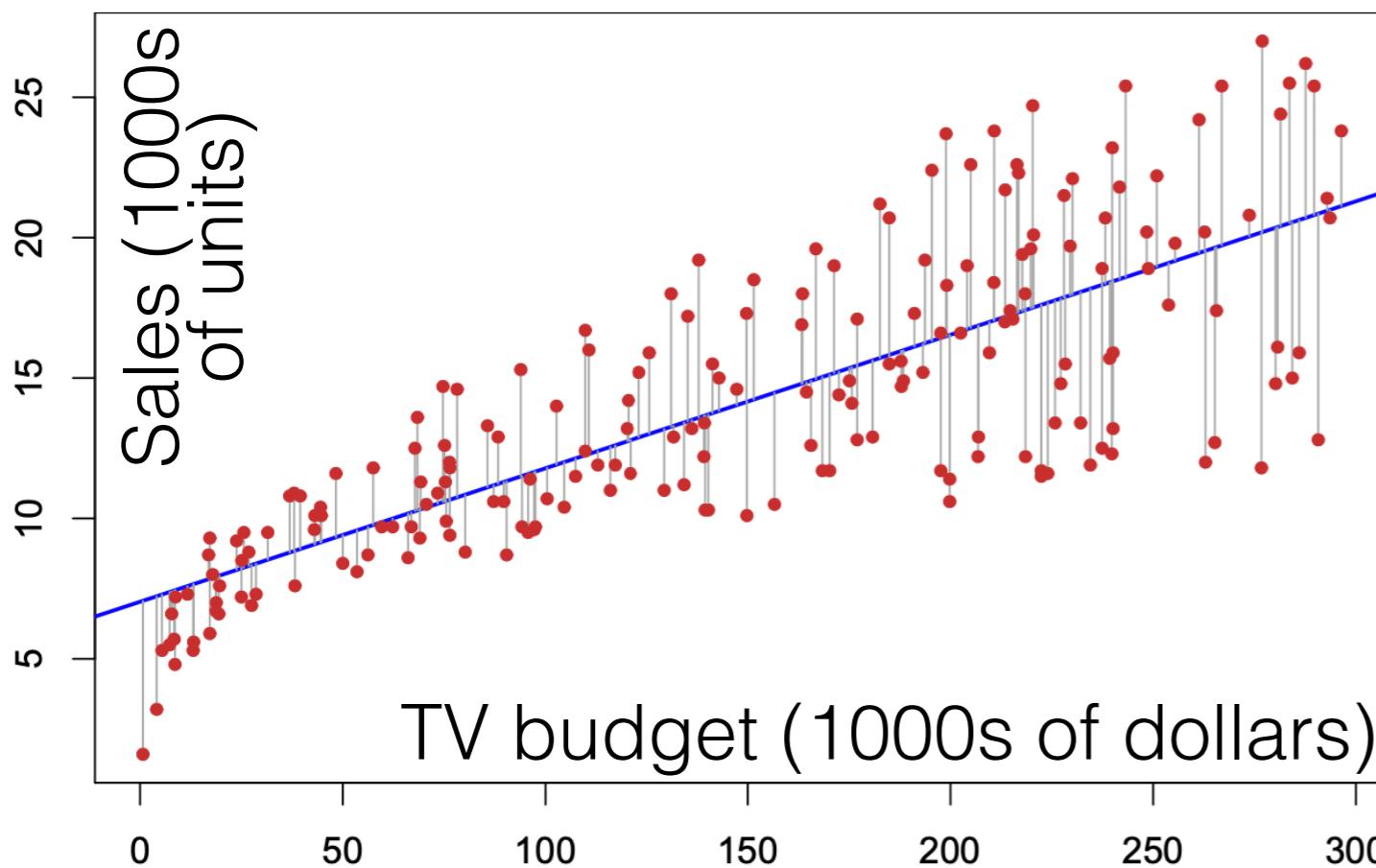


- Linear predictor and residuals for one feature (with an intercept)
- Can use residual plots of the fitted predictor to check model assumptions

[An Intro to Statistical Learning 2023,

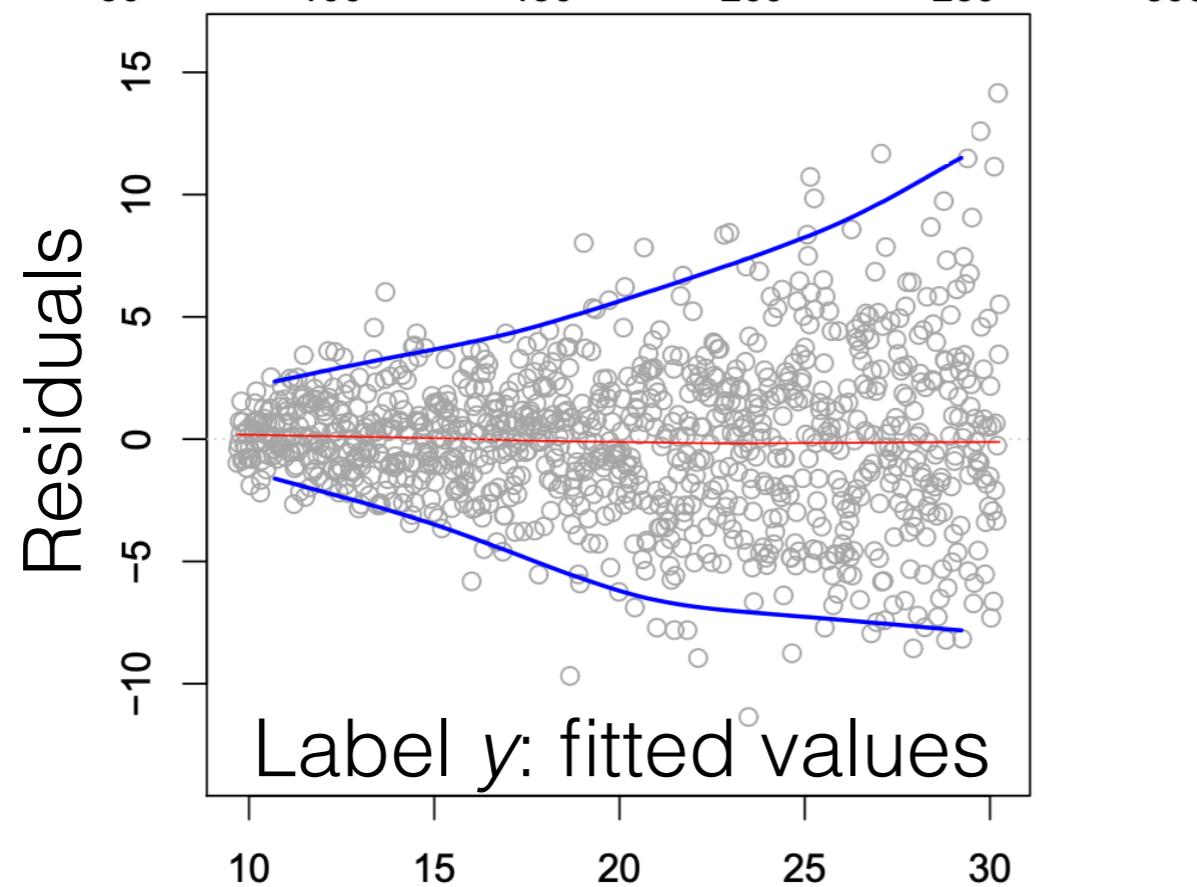
Fig 3.1]

Some visualizations

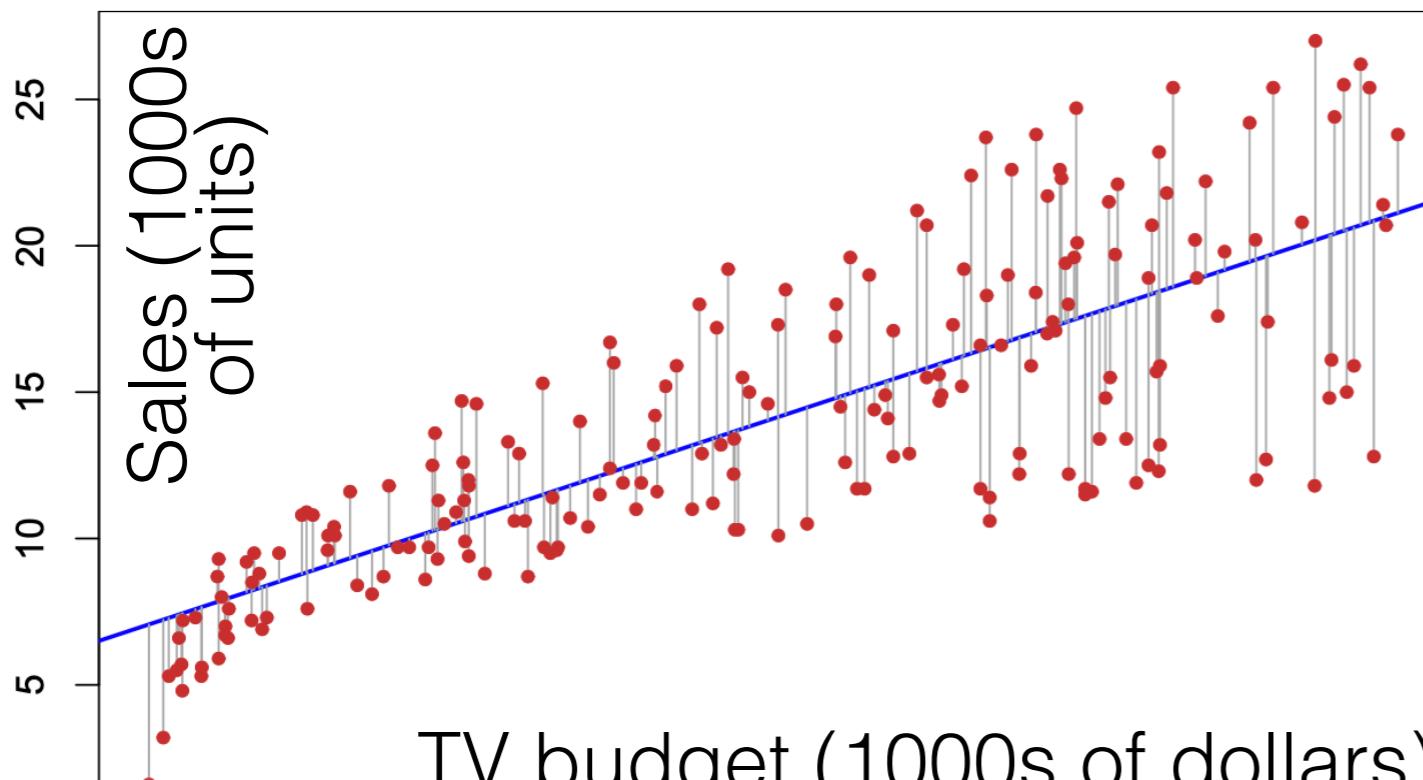


- Linear predictor and residuals for one feature (with an intercept)
- Can use residual plots of the fitted predictor to check model assumptions

[An Intro to Statistical Learning 2023,
Figs 3.1 and 3.11]

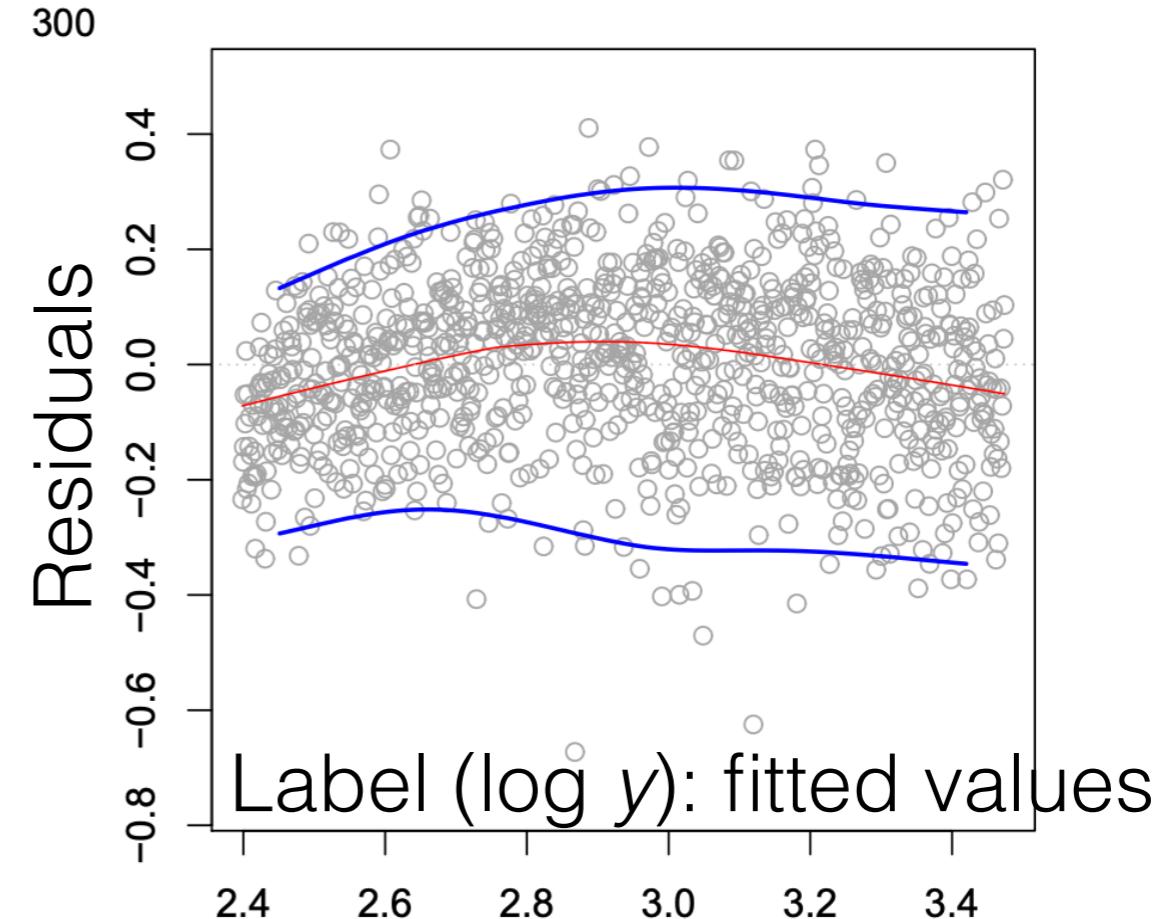
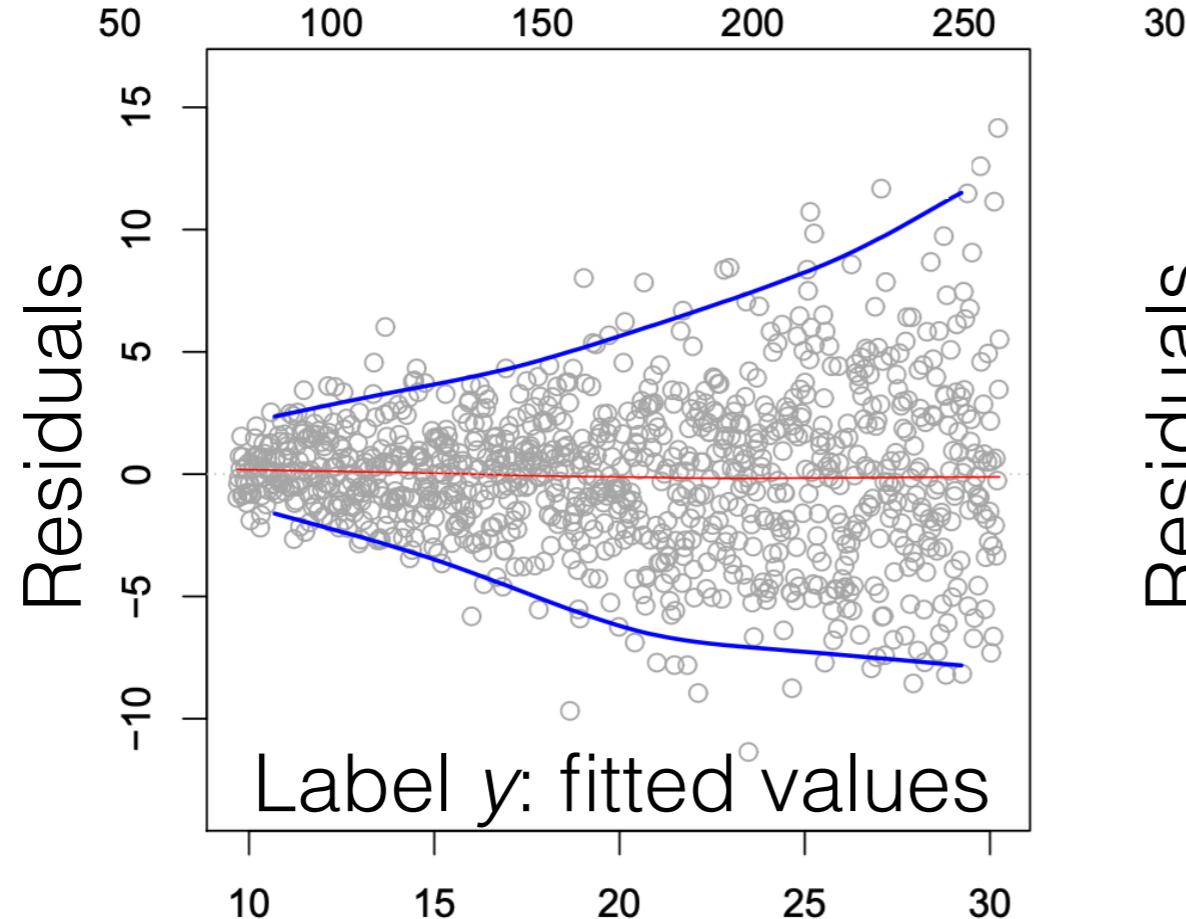


Some visualizations

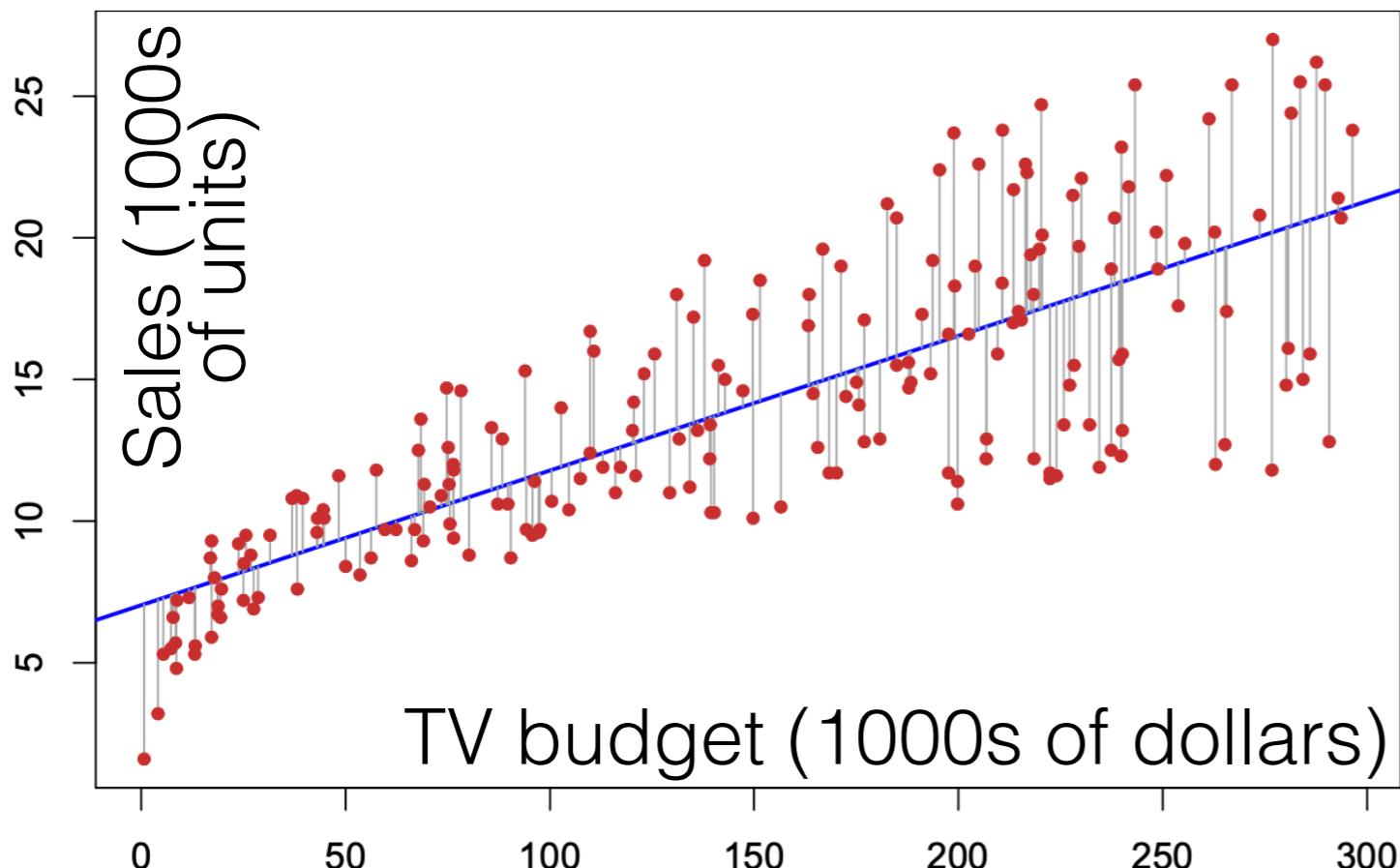


- Linear predictor and residuals for one feature (with an intercept)
- Can use residual plots of the fitted predictor to check model assumptions

[An Intro to Statistical Learning 2023,
Figs 3.1 and 3.11]



Some visualizations

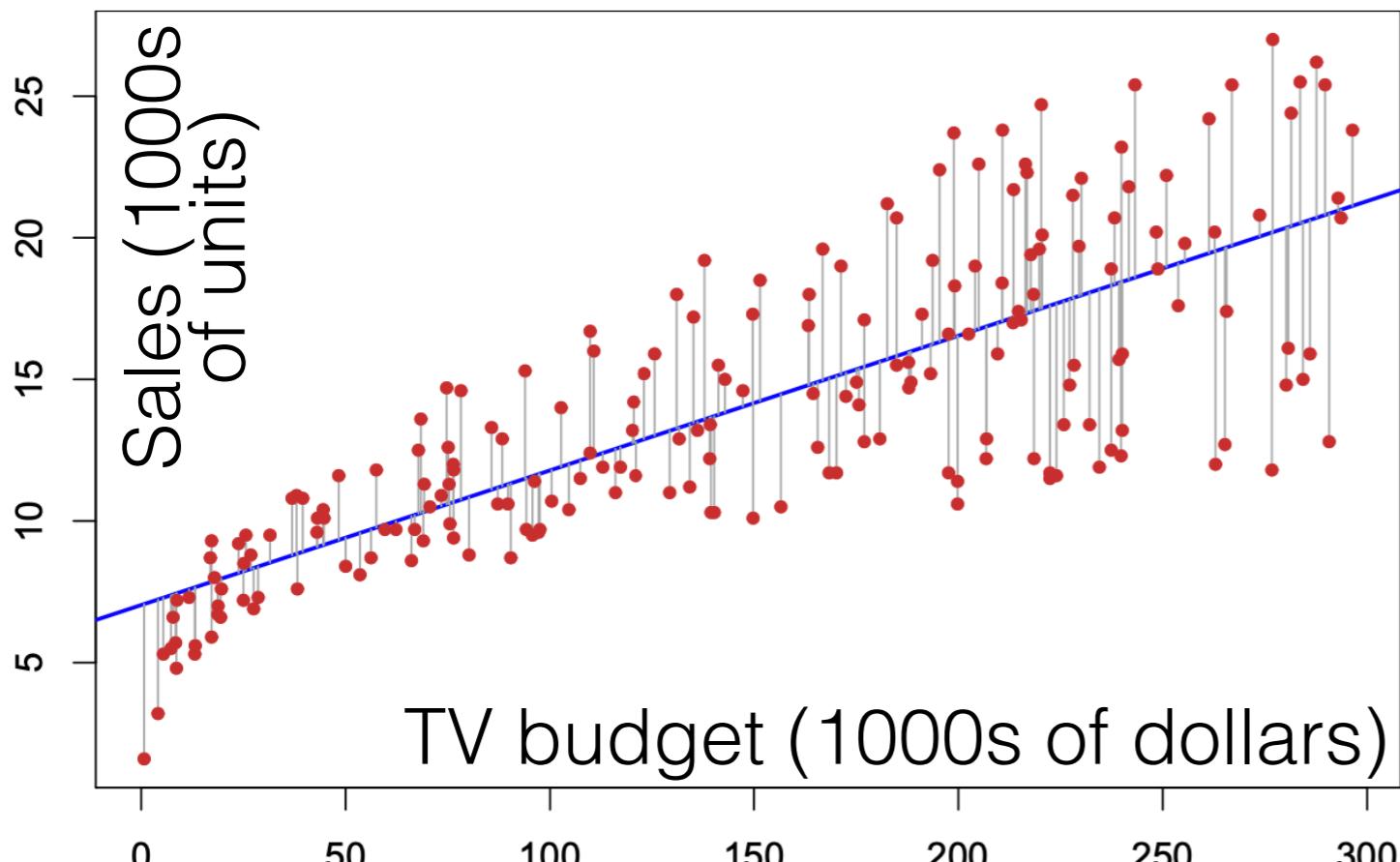


- Linear predictor and residuals for one feature (with an intercept)
- Can use residual plots of the fitted predictor to check model assumptions

[An Intro to Statistical Learning 2023,

Fig 3.1]

Some visualizations



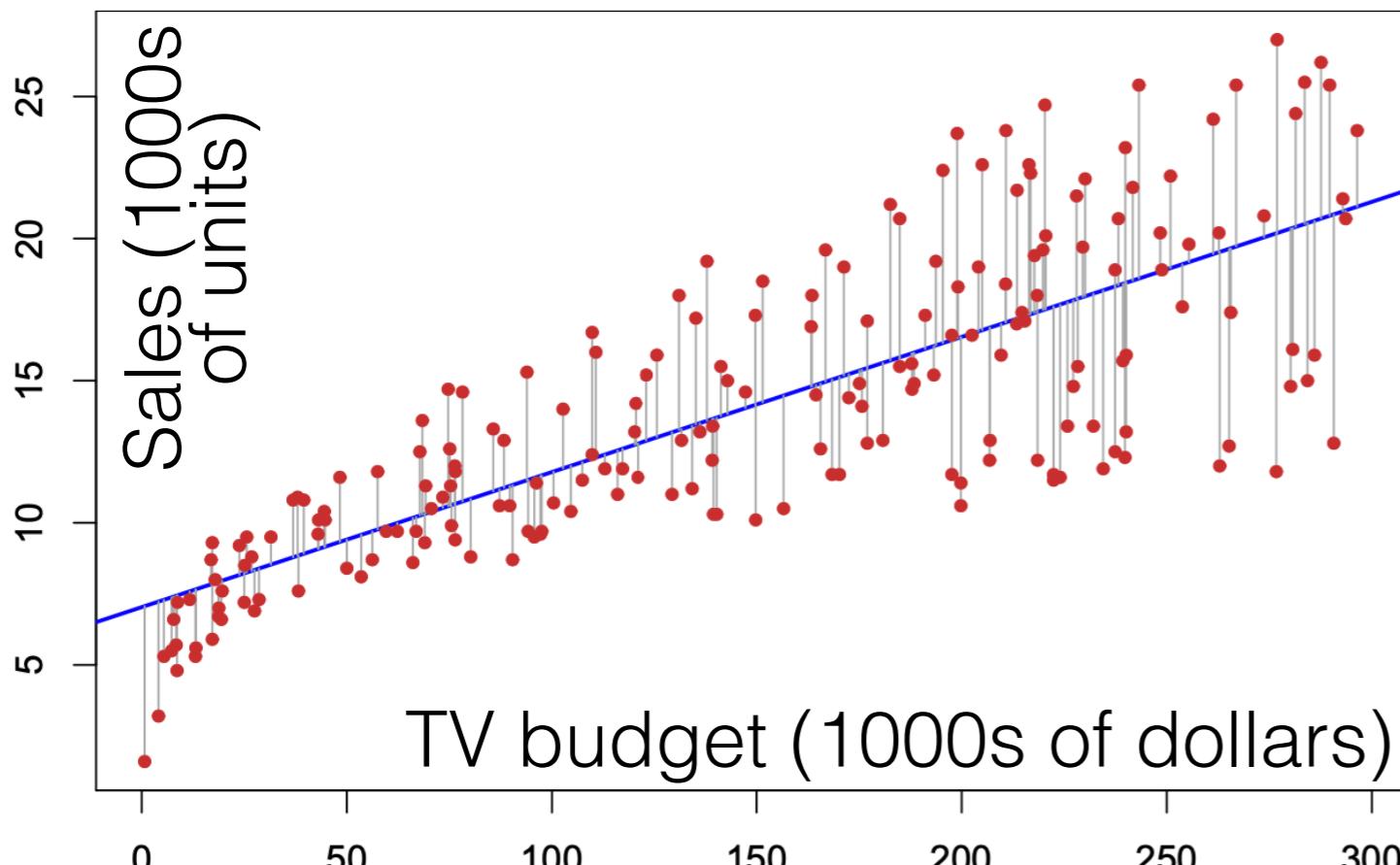
- Linear predictor and residuals for one feature (with an intercept)
- Can use residual plots of the fitted predictor to check model assumptions

[An Intro to Statistical Learning 2023,

Fig 3.1]

- Linear predictor & residuals for two features (with an intercept term)

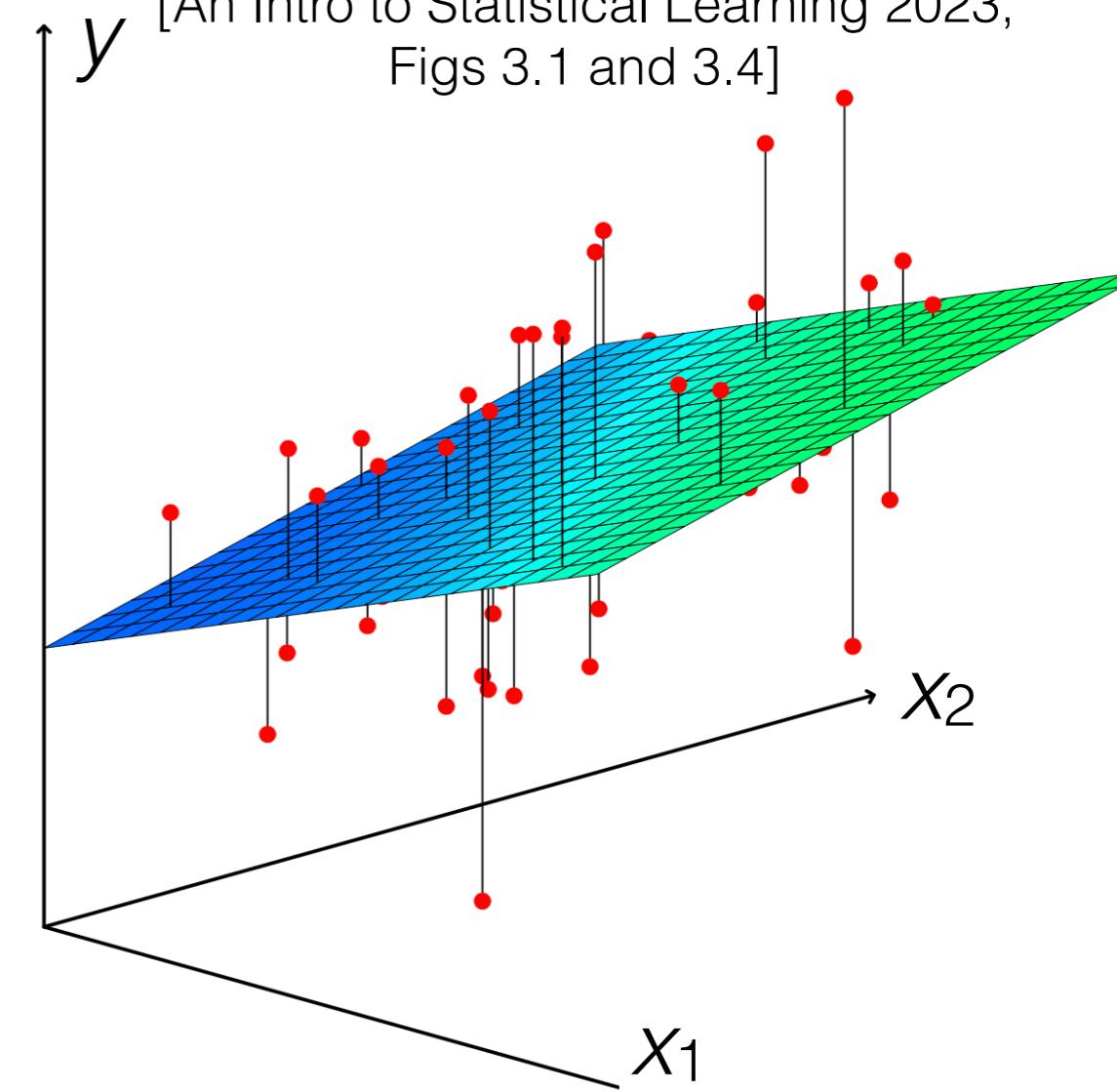
Some visualizations



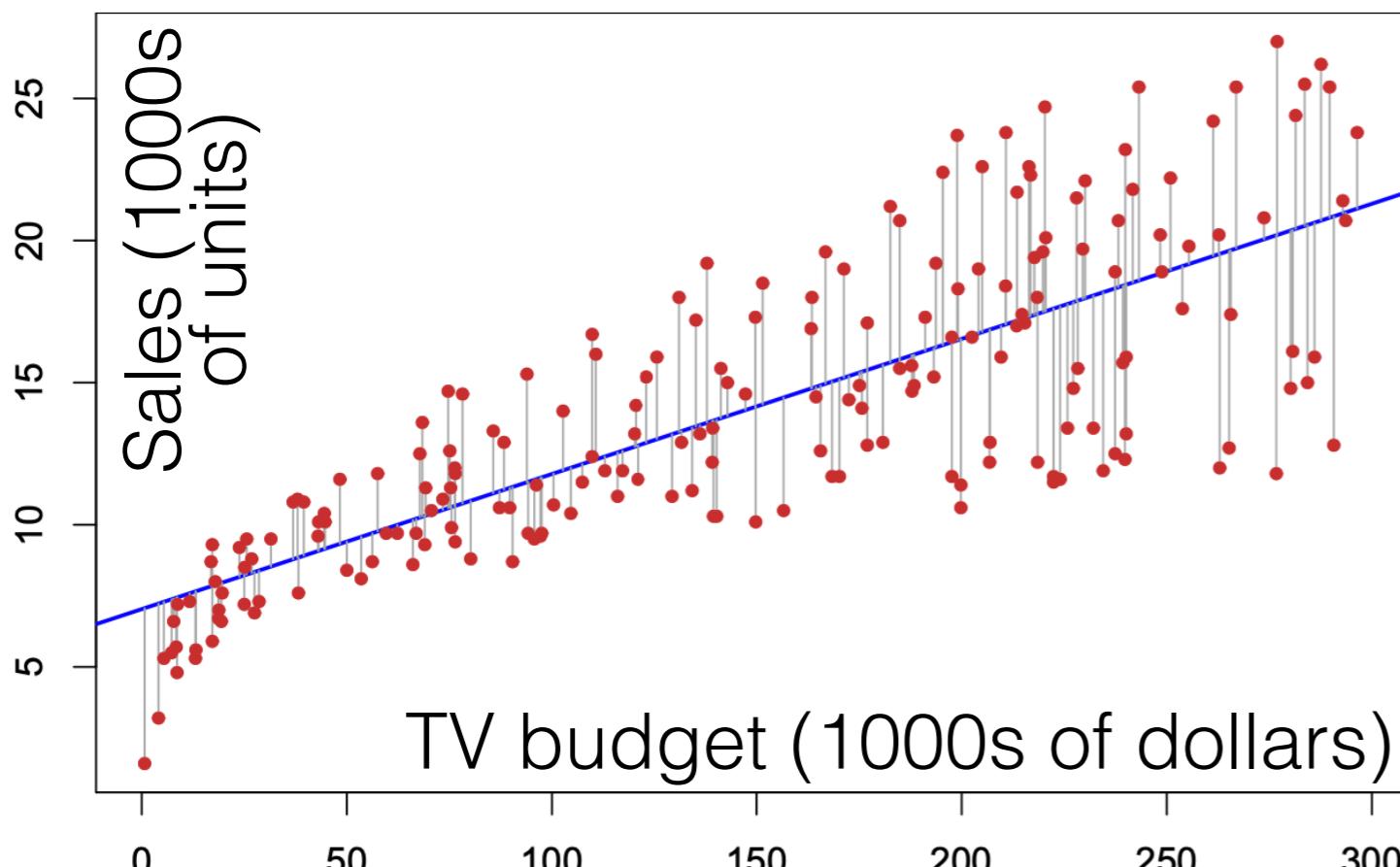
- Linear predictor & residuals for two features (with an intercept term)

- Linear predictor and residuals for one feature (with an intercept)
- Can use residual plots of the fitted predictor to check model assumptions

[An Intro to Statistical Learning 2023,
Figs 3.1 and 3.4]



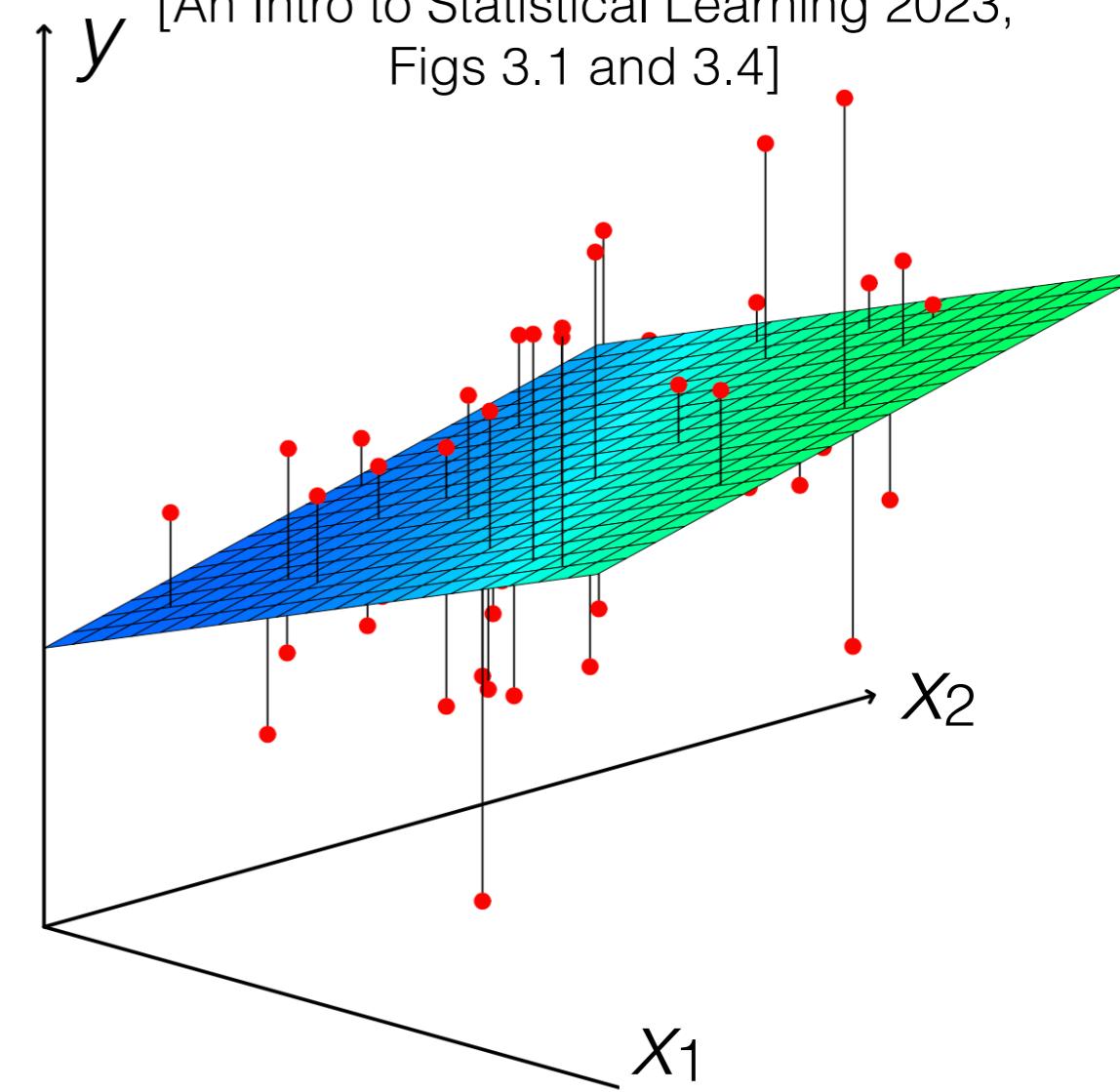
Some visualizations



- Linear predictor and residuals for one feature (with an intercept)
- Can use residual plots of the fitted predictor to check model assumptions

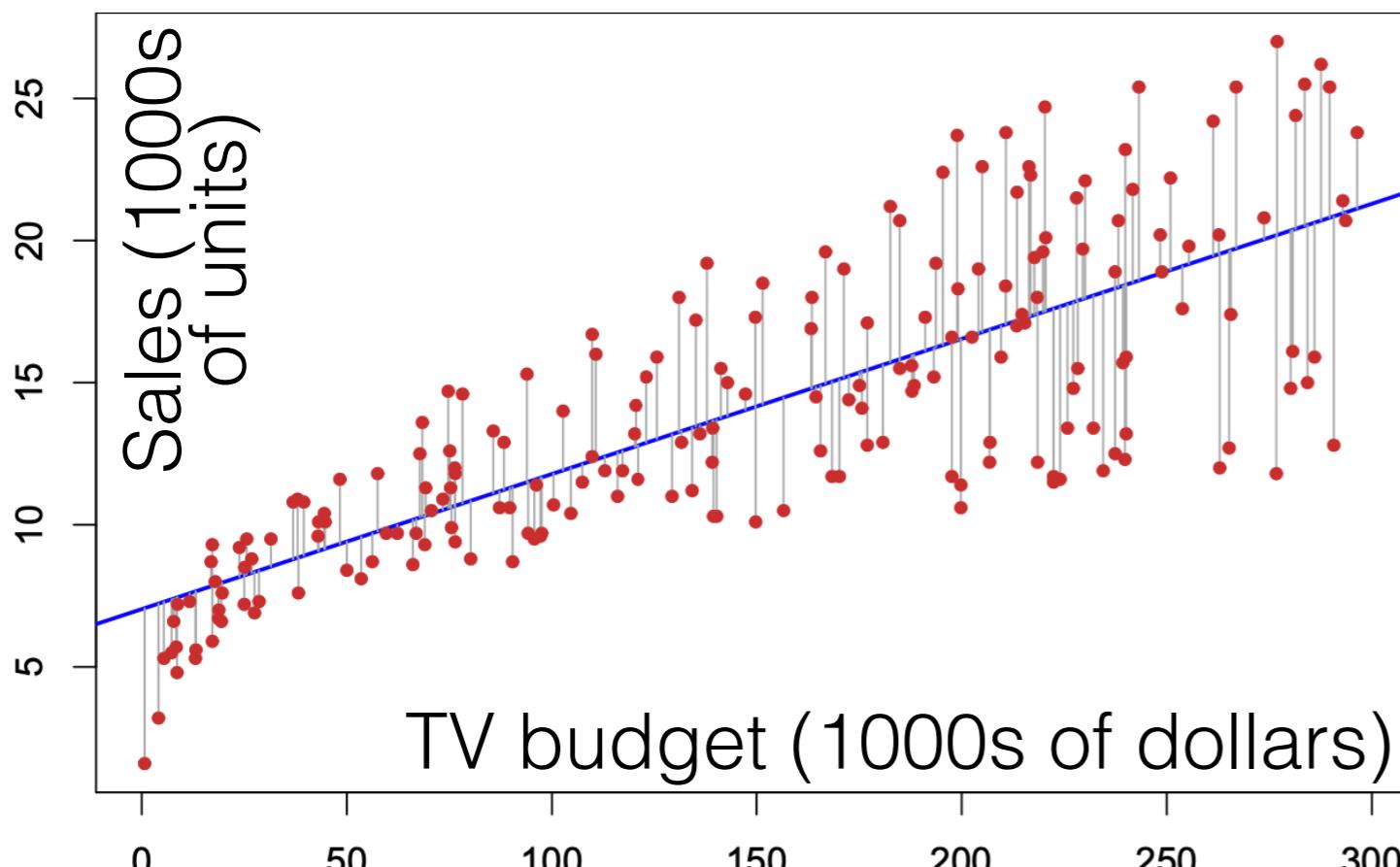
[An Intro to Statistical Learning 2023,

Figs 3.1 and 3.4]



- Linear predictor & residuals for two features (with an intercept term)
- In general, $\theta^\top x$ is a hyperplane (one dim less than (x,y) space)

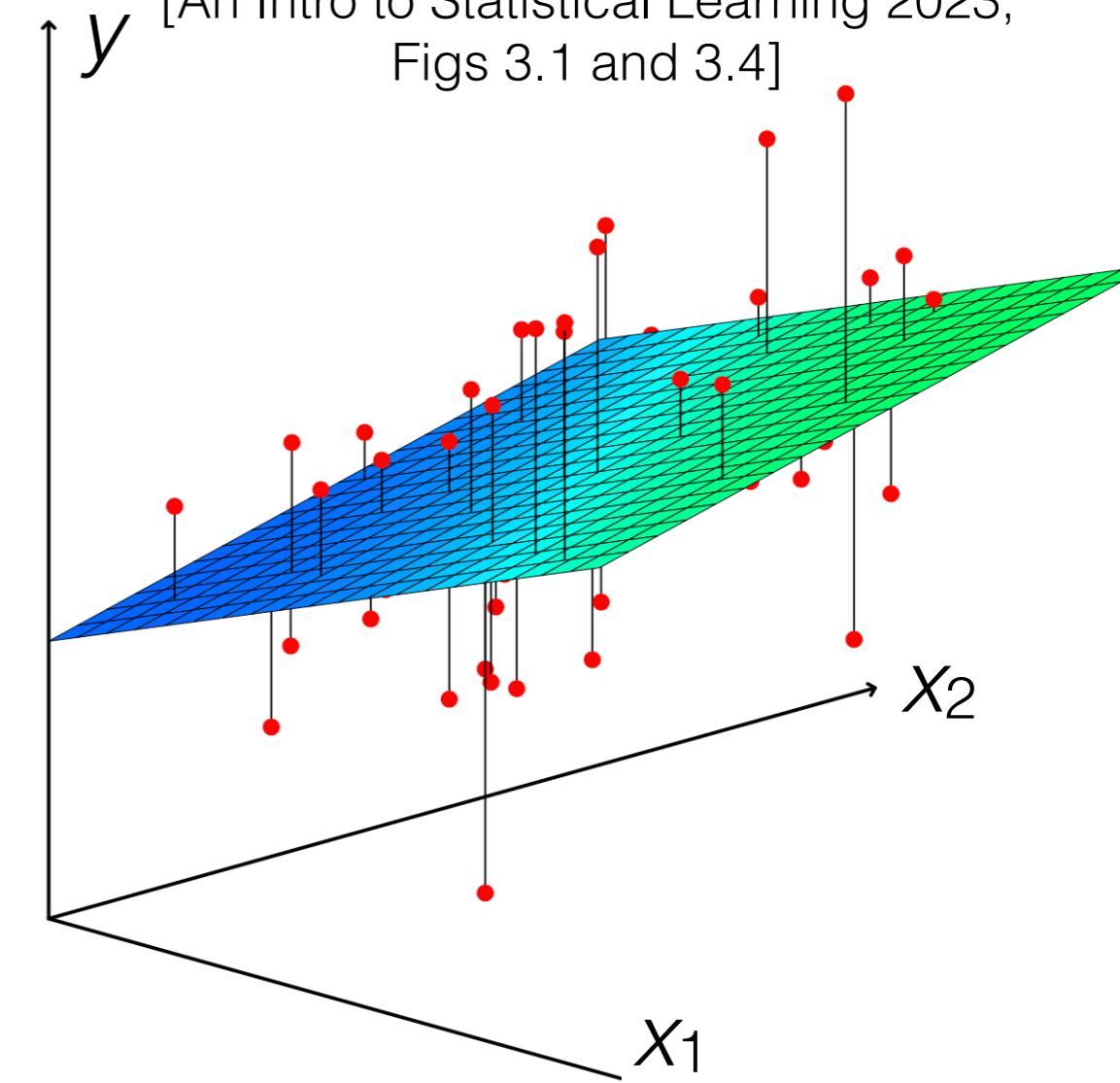
Some visualizations



- Linear predictor and residuals for one feature (with an intercept)
- Can use residual plots of the fitted predictor to check model assumptions

[An Intro to Statistical Learning 2023,

Figs 3.1 and 3.4]



- Linear predictor & residuals for two features (with an intercept term)
- In general, $\theta^\top x$ is a hyperplane (one dim less than (x,y) space)
- Check: regressing label on the features is not equivalent to regressing one feature on label and remaining features

Empirical risk minimization

Empirical risk minimization

- Empirical risk for square loss:

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) =$$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
 - Then minimizing the empirical risk minimizes

$$\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2$$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2 \quad \text{same objective from maximum likelihood! (up to a constant factor)}$$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2$$

same objective from maximum likelihood! (up to a constant factor)

- relation is specific to this model

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- Notation:

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

- Notation:

$$X = \begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix}$$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{array}{l} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{array}$$

- relation is specific to this model

- Notation:

$$X = \begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix}$$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

- Notation:

$$X = \begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix}$$

NxD

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

- Notation:

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(N)}]^\top$$

$$X = \begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix}$$

NxD

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

- Notation:

$$X = \begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix} \quad \begin{matrix} \text{Nx1} \\ \text{Dx1} \end{matrix}$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(N)}]^\top$$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

- Notation:

$$X = \begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix} \quad \begin{matrix} \text{Nx1} \\ \text{DxD} \end{matrix}$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(N)}]^\top$$

- Then $\text{RSS}(\theta) =$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

- Notation:

$$X = \begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix} \quad \begin{matrix} \text{Nx1} \\ \text{DxD} \end{matrix}$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(N)}]^\top$$

- Then $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

- Notation:

$$X = \begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix} \quad \begin{matrix} \text{Nx1} \\ \text{Dx1} \end{matrix}$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(N)}]^\top$$

- Then $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$

$$\begin{matrix} \text{NxD} \\ \text{Nx1} \end{matrix}$$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

- Notation:

$$X = \begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix} \quad \begin{matrix} \text{Nx1} \\ \text{Dx1} \end{matrix}$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(N)}]^\top$$

- Then $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$

$\begin{matrix} \text{NxD} & \text{Dx1} & \text{Nx1} \end{matrix}$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

- Notation:

$$X = \begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix} \quad \begin{matrix} \text{Nx1} \\ \text{Dx1} \end{matrix}$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(N)}]^\top$$

- Then $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$

$$\begin{matrix} \text{Nx1} & \text{Dx1} & \text{Nx1} \\ & & \overbrace{\quad\quad\quad}^{\text{Nx1}} \end{matrix}$$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

- Notation:

$$X = \underbrace{\begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix}}_{NxD}$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(N)}]^\top$$

Nx1

- Then $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$

NxD Dx1 Nx1

$\underbrace{}_{1xN} \underbrace{}_{Nx1}$

Empirical risk minimization

- Empirical risk for square loss:

$$\frac{1}{N} \sum_{n=1}^N L(y^{(n)}, h(x^{(n)})) = \frac{1}{N} \sum_{n=1}^N (y^{(n)} - h(x^{(n)}))^2$$

- Still not useful if we allow all decision rules (functions) h
- What if we instead allow only linear predictors: $h(x) = \theta^\top x$
- Then minimizing the empirical risk minimizes

$$\underbrace{\frac{1}{N} \sum_{n=1}^N (y^{(n)} - \theta^\top x^{(n)})^2}_{\text{Mean squared error (MSE)}} \quad \begin{matrix} \text{same objective from maximum} \\ \text{likelihood! (up to a constant factor)} \end{matrix}$$

- relation is specific to this model

- Notation:

$$X = \underbrace{\begin{pmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(N)})^\top \end{pmatrix}}_{NxD}$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(N)}]^\top$$

Nx1

- Then $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$

$$\underbrace{(X\theta - Y)^\top}_{\substack{\text{NxD} \\ 1xN}} \underbrace{(X\theta - Y)}_{\substack{\text{Dx1} \\ 1x1}}$$

$\underbrace{\qquad\qquad\qquad}_{\substack{\text{Nx1} \\ 1x1}}$

Let's optimize!

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the Dx1 vector with d th element: $\partial f / \partial \theta_d$

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the Dx1 vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the DxD matrix with (d,d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the Dx1 vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the DxD matrix with (d,d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions:

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the Dx1 vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the DxD matrix with (d,d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) =$

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the Dx1 vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the DxD matrix with (d,d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top (X\theta - Y)$

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the Dx1 vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the DxD matrix with (d,d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top (X\theta - Y)$

check you can derive! 

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the Dx1 vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the DxD matrix with (d,d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top(X\theta - Y) \stackrel{\text{set}}{=} 0$

check you can derive! 

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the Dx1 vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the DxD matrix with (d,d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top (X\theta - Y) \stackrel{\text{set}}{=} 0$

check you can derive!  Dx1

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the Dx1 vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the DxD matrix with (d,d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2\underbrace{X^\top}_{\text{Dx1}}(X\theta - Y) \stackrel{\text{set}}{=} 0$
- check you can derive! 

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2\underbrace{X^\top}_{D \times 1} \underbrace{(X\theta - Y)}_{D \times N \quad Nx1} \stackrel{\text{set}}{=} 0$
- check you can derive! 

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2\underbrace{X^\top}_{D \times 1} \underbrace{(X\theta - Y)}_{D \times N} \stackrel{\text{set}}{=} 0$
- check you can derive! 

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top \underbrace{(X\theta - Y)}_{\substack{\text{Dx1} \\ \text{DxN} \quad \text{Nx1}}} = 0$
check you can derive! 
 - Second-order conditions:

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top \underbrace{(X\theta - Y)}_{\substack{\text{Dx1} \\ \text{DxN} \\ \text{Nx1}}} = 0$
check you can derive!  set
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) =$

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top (X\theta - Y) \xrightarrow{\text{set}} 0$
check you can derive! 
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2X^\top X$

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2 \underbrace{X^\top}_{D \times 1} \underbrace{(X\theta - Y)}_{D \times N} \underbrace{\nabla \text{RSS}(\theta)}_{N \times 1} = 0$
check you can derive! 
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2 \underbrace{X^\top X}_{D \times D}$

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2 \underbrace{X^\top}_{D \times 1} \underbrace{(X\theta - Y)}_{D \times N} \underbrace{\nabla \text{RSS}(\theta)}_{N \times 1} = 0$ set
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2 \underbrace{X^\top X}_{D \times D}$

check you can derive!

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2\underbrace{X^\top}_{D \times 1} \underbrace{(X\theta - Y)}_{D \times N} \stackrel{\text{set}}{=} 0$
check you can derive! 
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2\underbrace{X^\top X}_{D \times D} \stackrel{\text{set}}{>} 0$ is positive definite

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top(X\theta - Y) \stackrel{\text{set}}{=} 0$
check you can derive! 
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2X^\top X \stackrel{\text{set}}{>} 0$
For now, suppose $N > D$ & X is full rank 
is positive definite

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top (X\theta - Y) \stackrel{\text{set}}{=} 0$
check you can derive! 
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2X^\top X \stackrel{\text{set}}{>} 0$
 - For now, suppose $N > D$ & X is full rank
 - I.e. X has rank D , i.e. columns linearly independent

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top (X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top (X\theta - Y) \stackrel{\text{set}}{=} 0$
check you can derive! 
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2X^\top X \stackrel{\text{set}}{>} 0$
 - For now, suppose $N > D$ & X is full rank
 - I.e. X has rank D , i.e. columns linearly independent
 - Then $X^\top X$ is positive definite and invertible

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top(X\theta - Y) \stackrel{\text{set}}{=} 0$
check you can derive! 
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2X^\top X \stackrel{\text{set}}{>} 0$
 - For now, suppose $N > D$ & X is full rank
 - I.e. X has rank D , i.e. columns linearly independent
 - Then $X^\top X$ is positive definite and invertible
 - And we can solve the first-order conditions:

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top(X\theta - Y) \stackrel{\text{set}}{=} 0$
check you can derive! 
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2X^\top X \stackrel{\text{set}}{>} 0$
 - For now, suppose $N > D$ & X is full rank
 - I.e. X has rank D , i.e. columns linearly independent
 - Then $X^\top X$ is positive definite and invertible
 - And we can solve the first-order conditions:

$$X^\top X \hat{\theta} = X^\top Y$$

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the Dx1 vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the DxD matrix with (d,d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top(X\theta - Y) \stackrel{\text{set}}{=} 0$
check you can derive! 
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2X^\top X \stackrel{\text{set}}{>} 0$
 - For now, suppose $N > D$ & X is full rank
 - I.e. X has rank D , i.e. columns linearly independent
 - Then $X^\top X$ is positive definite and invertible
 - And we can solve the first-order conditions:

$$\begin{aligned} X^\top X \hat{\theta} &= X^\top Y \\ \hat{\theta} &= (X^\top X)^{-1} X^\top Y \end{aligned}$$

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the Dx1 vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the DxD matrix with (d,d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top(X\theta - Y) \stackrel{\text{set}}{=} 0$
check you can derive! 
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2X^\top X \stackrel{\text{set}}{>} 0$
 - For now, suppose $N > D$ & X is full rank
 - I.e. X has rank D , i.e. columns linearly independent
 - Then $X^\top X$ is positive definite and invertible
 - And we can solve the first-order conditions:

$$\begin{aligned} X^\top X \hat{\theta} &= X^\top Y \\ \hat{\theta} &= (X^\top X)^{-1} X^\top Y \end{aligned}$$

“ordinary least squares” (OLS)

Let's optimize!

- We want to minimize $\text{RSS}(\theta) = (X\theta - Y)^\top(X\theta - Y)$
 - Notation: for function $f(\theta)$, the gradient $\nabla_\theta f$ is the $D \times 1$ vector with d th element: $\partial f / \partial \theta_d$
 - The Hessian $\nabla_\theta^2 f$ is the $D \times D$ matrix with (d, d') element $\frac{\partial^2 f}{\partial \theta_d \partial \theta_{d'}}$
 - First-order conditions: $\nabla \text{RSS}(\theta) = 2X^\top(X\theta - Y) \stackrel{\text{set}}{=} 0$
check you can derive! 
 - Second-order conditions: $\nabla^2 \text{RSS}(\theta) = 2X^\top X \stackrel{\text{set}}{>} 0$
 - For now, suppose $N > D$ & X is full rank
 - I.e. X has rank D , i.e. columns linearly independent
 - Then $X^\top X$ is positive definite and invertible
 - And we can solve the first-order conditions:

$$\begin{aligned} X^\top X \hat{\theta} &= X^\top Y \\ \hat{\theta} &= (X^\top X)^{-1} X^\top Y \end{aligned}$$

“ordinary least squares” (OLS)

- A closed-form solution isn't always better
 - Matrix inversion can be expensive for large dimensions

References (1/1)

Castro Torres, Andrés F., and Aliakbar Akbaritabar. "The use of linear models in quantitative research." *Quantitative Science Studies* (2024): 1-21.

Chen, Danqi, Jason Bolton, and Christopher D. Manning. "A thorough examination of the CNN/Daily Mail reading comprehension task." *ACL* (2016).

Lipton, Zachary C., and Jacob Steinhardt. "Troubling Trends in Machine Learning Scholarship: Some ML papers suffer from flaws that could mislead the public and stymie future research." *Queue* 17.1 (2019): 45-77.

Salganik, Matthew J., et al. "Measuring the predictability of life outcomes with a scientific mass collaboration." *Proceedings of the National Academy of Sciences* 117.15 (2020): 8398-8403.