

Linear Regression

SWC Neuroinformatics 2024

Lior Fox
Gatsby Unit, UCL

Intro

Linear regression is ***everywhere***

The practicalities are endless

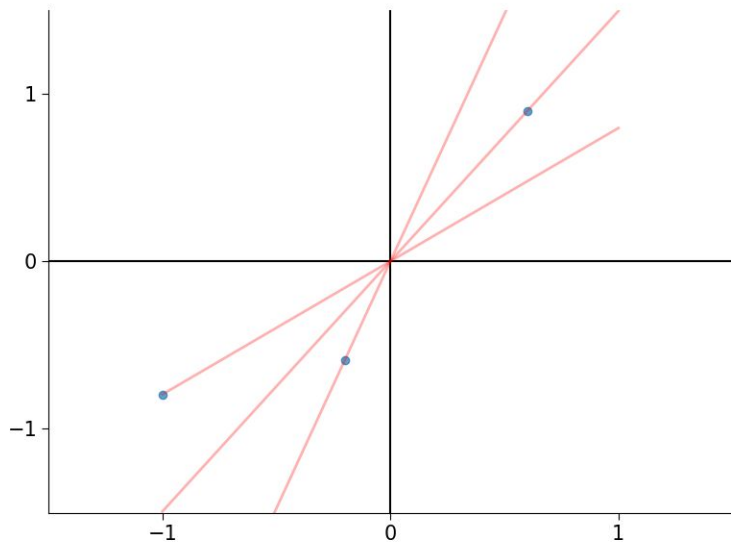
It is worthwhile taking time to go through the basics

Outline

- The Least Squares solution in 1 and multiple dimensions
- Model complexity, Bias-Variance tradeoff, and regularization
- Probabilistic / Bayesian interpretation of linear regression
- Some examples

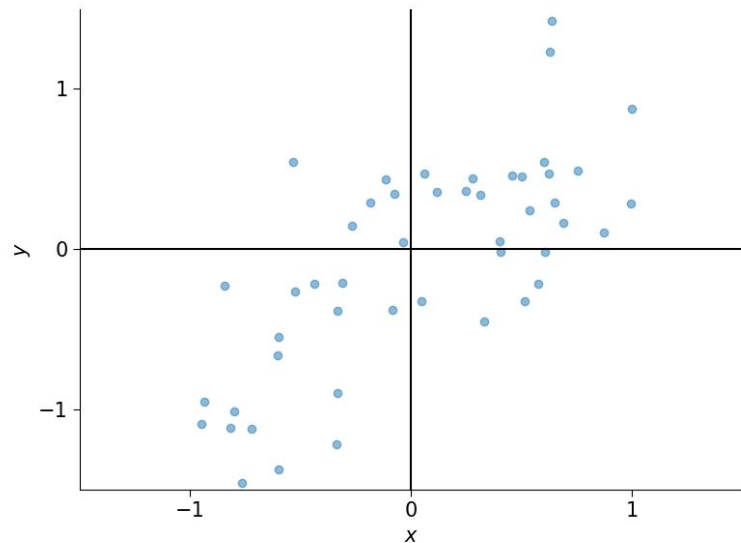
The problem

- Given a sample (x^i, y^i) , find a good way to describe the dependence of y on x .
- Immediate questions:
 - What makes a fit “good”?
 - What type of dependences are we willing to consider?
- Keep in mind: we would like our model to be able to predict *unseen* future data
- Restrict to **linear** model: $y \sim bx$
 - If there was a single point, we could fit perfectly
 - What shall we do if there are multiple points?



Least Squares

- Instead of trying to perfectly fit each point, minimise the sum of squared errors
- The minimisation is over the parameter, b
- Why squared error?
 - Mathematically convenient
 - The “right” thing under some assumptions (later)
- Deriving the solution



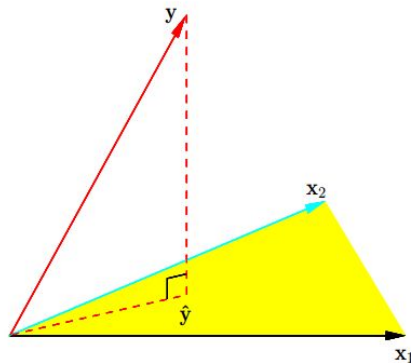
Multiple predictors

- Each \mathbf{x} is now a vector $[x_1 \ x_2 \ \dots \ x_p]^T$
- Linear mapping $\mathbf{x} \mapsto \mathbf{y}$ is parameterised by a vector of coefficients, $\mathbf{y} \sim \mathbf{x}^T \mathbf{b}$

Deriving the solution

Interpretation:

- The analogy to the 1D case
- The predictions for \mathbf{y} are $\mathbf{Xb} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$: the *orthogonal projection* of \mathbf{y} onto the space spanned by the input vectors



Adding an offset

- By positing $y \sim \mathbf{b}x$ we have constrained the regression line to pass through the origin
- We can instead assume $y \sim \mathbf{b}x + \mathbf{c}$ to allow an offset (aka ‘bias’, ‘intercept’)

$$\mathbf{b}^* = \text{Cov}[x, y] / \text{Var}[x]$$

$$\mathbf{c}^* = \mathbf{E}[y] - \mathbf{b}^* \mathbf{E}[x]$$

- The optimal (minimising MSE) line goes through the sample average
- Therefore we can alternatively *center* the data, and consider the homogeneous model
 - Note that this recovers the “full” solution, as $\text{Cov}[x, y] = \mathbf{E}[xy]$ and $\text{Var}[x] = \mathbf{E}[x^2]$
- Alternatively, we can view $y \sim \mathbf{b}x + \mathbf{c}$ as a multiple predictors case
- For that, we redefine the examples, with mapping $x \mapsto [x, 1]$.
- Exercise: show the solutions match: $[\mathbf{b}^* \ \mathbf{c}^*]^T = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$, after the remapping

The asymmetry in regression

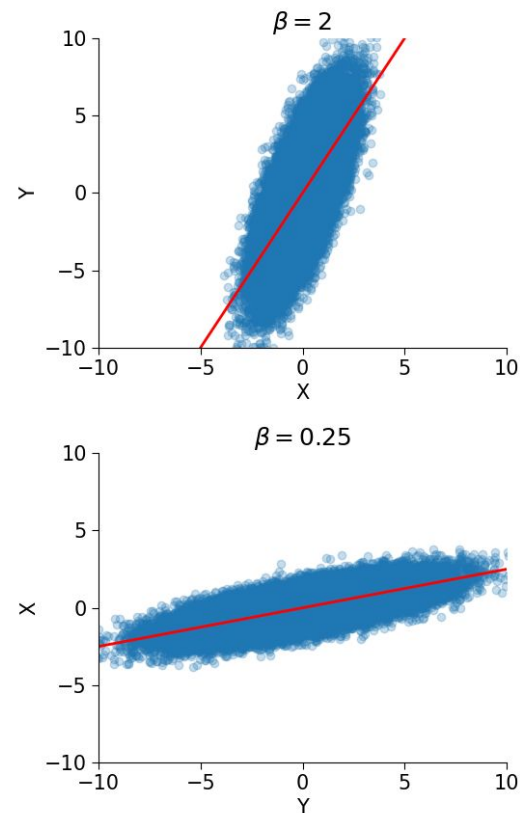
- Regressing $y \sim \mathbf{b}x$ and $x \sim \mathbf{a}y$ isn't the same – in general, $a \neq 1/b$
- This is due to the division by $\text{Var}[x]$!

$$b_{y \sim x} = E[xy]/E[x^2]$$

$$b_{x \sim y} = E[xy]/E[y^2]$$

- Changes the role of the “noise”

In this example: $x \sim N(0,1)$, $y=2x+e$, $e \sim N(0,2^2)$



Inference and hypothesis testing

- So far we didn't commit too much to the true data distribution
- If we add some assumptions, we can say more about the result

In particular, we assume that:

- observations \mathbf{y} were generated by an **unknown** linear model
- additive gaussian noise
- independent (and equal variance) noise between different observations

Mathematically, assume that $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, with $\boldsymbol{\varepsilon} \sim \mathbf{N}(\mathbf{0}, \sigma^2 \mathbf{I})$

Inference and hypothesis testing

We can now work out the distribution of the estimator β – this is a simple exercise in manipulating gaussian distributions:

$$\begin{aligned}\beta &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \beta + \epsilon) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \\ &= \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

And therefore, we have $\beta \sim \mathbf{N}(\beta, (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2)$

- Unbiased
- Lowest variance among all linear unbiased estimators (*Gauss Markov Theorem*)

Inference and hypothesis testing

- Now that we have the distribution of β we can perform hypothesis testing
- For example, we might want to test whether $\beta_i \neq 0$:
 - Write the distribution of β_i under the null, i.e., under $\beta_i = 0$, and compare the observed value
 - In our case, this would be a simple Z-test (since we assumed known variance)
 - If the variance is unknown, it can be estimated, resulting in a t -distribution for the standardised β_i
- The assumptions allow the usage of many standard statistical tools (t-test, F-test, etc)
- Can be slightly relaxed if we are willing to use nonparametric tests (e.g., bootstrap)

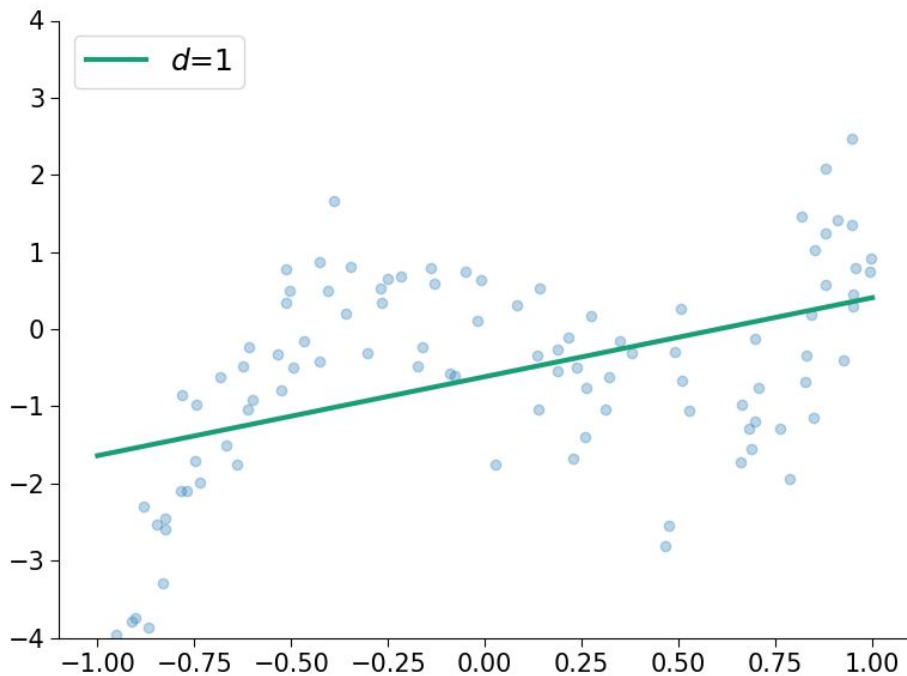
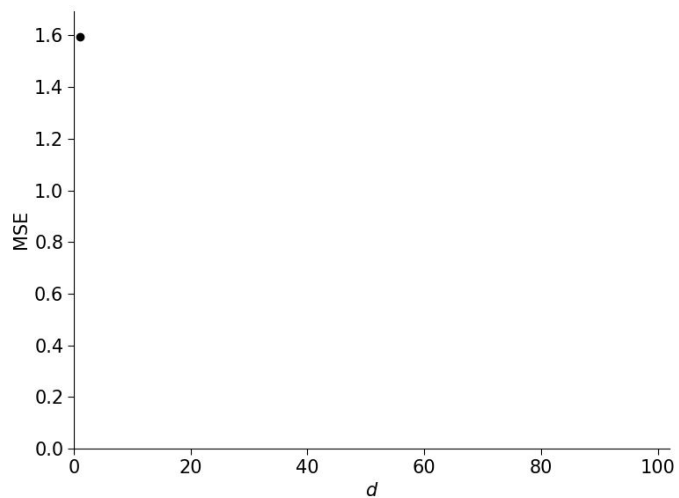
Rethinking the “linear” in linear regression

- The important thing was linearity **in the parameters**
- We can add more predictors, which are functions of the original variables, resulting in non-linear functions of those

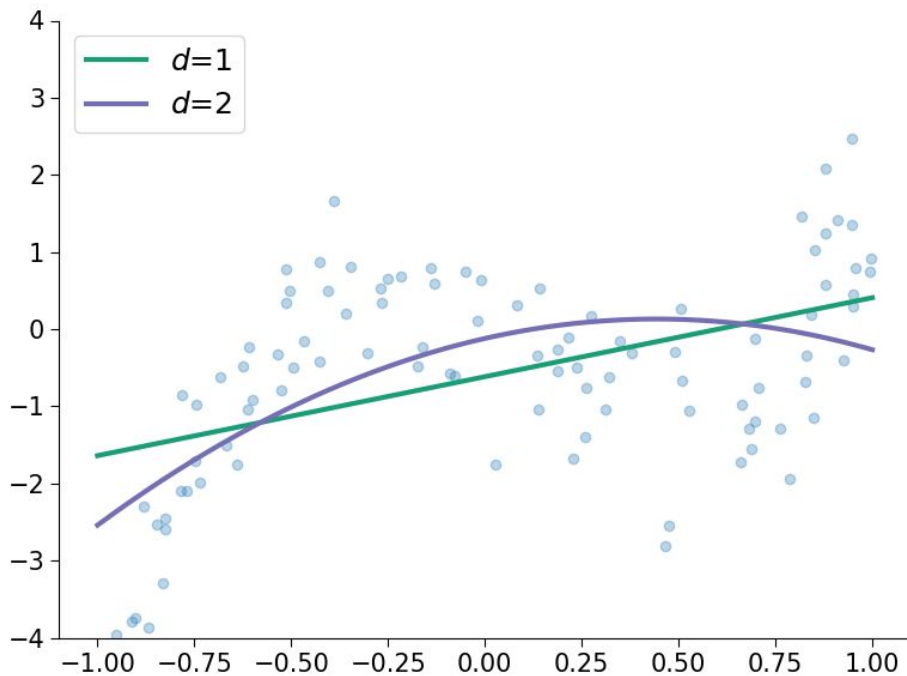
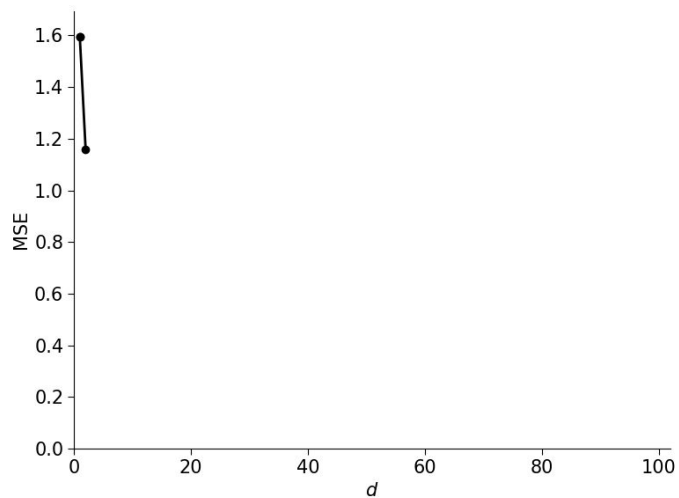
Example: polynomial regression

- Suppose we have 1D inputs x . We map $\mathbf{x} \mapsto [\mathbf{x} \ x^2 \ x^3 \ \dots \ x^d]^T$, and solve for \mathbf{b}
- Then effectively we have a model $\mathbf{y} \sim \mathbf{b}_1 \mathbf{x} + \mathbf{b}_2 \mathbf{x}^2 + \dots + \mathbf{b}_d \mathbf{x}^d$ – a d degree polynomial
- What happens as we increase \mathbf{d} ?
 - We are making the model less and less constrained
 - Since $\text{poly}(d) \subset \text{poly}(d+1)$, we can only decrease the error by taking larger and larger \mathbf{d}
 - Is this the right thing to do?

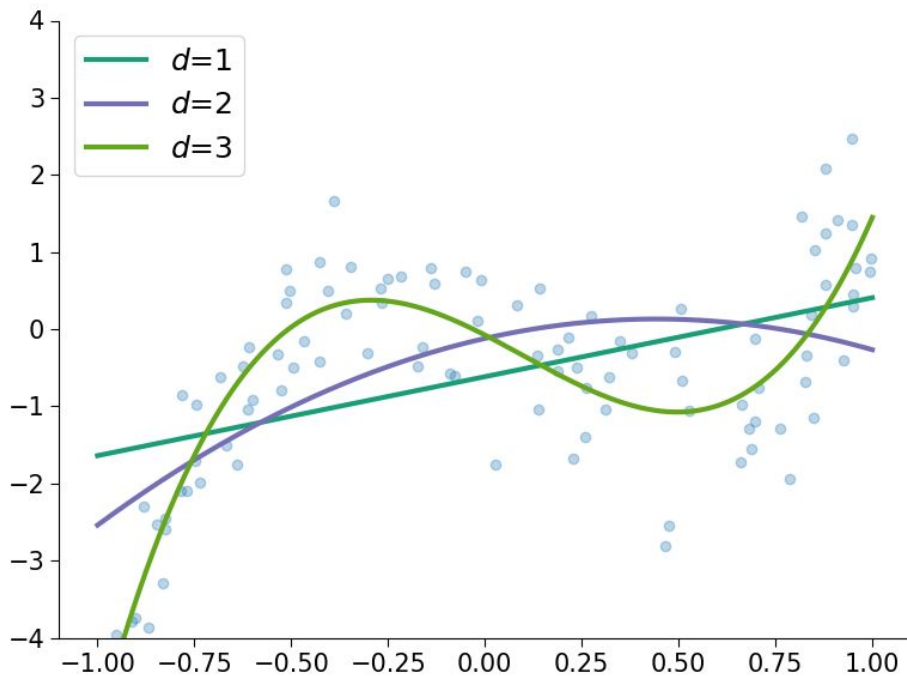
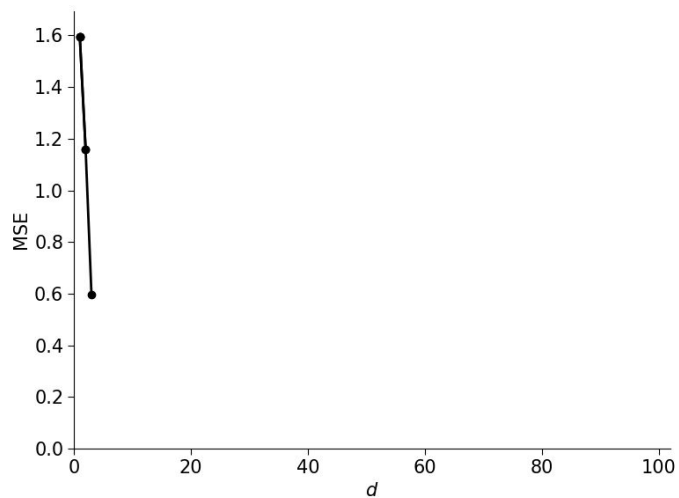
Model complexity



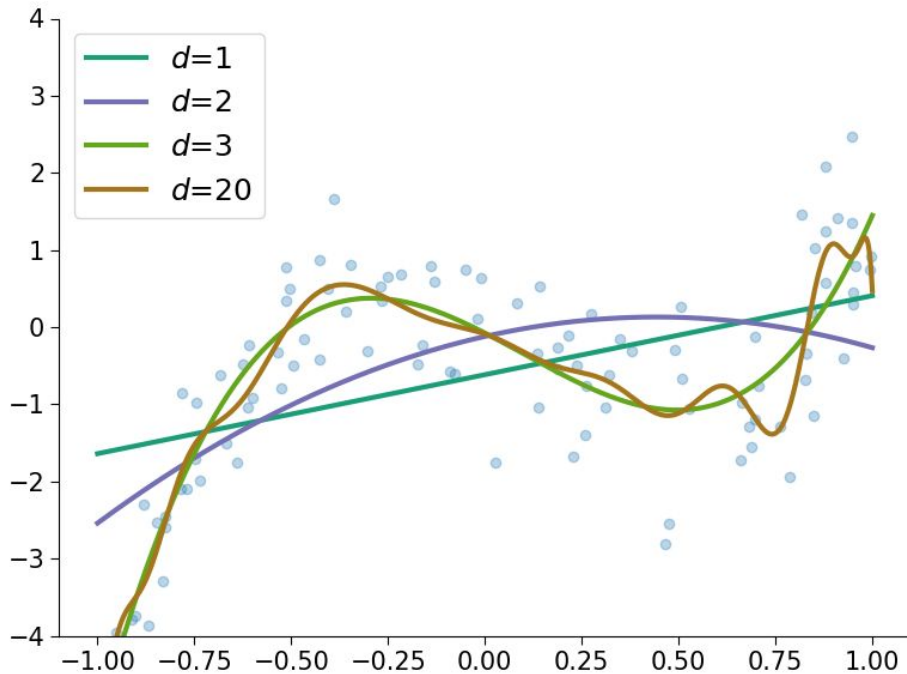
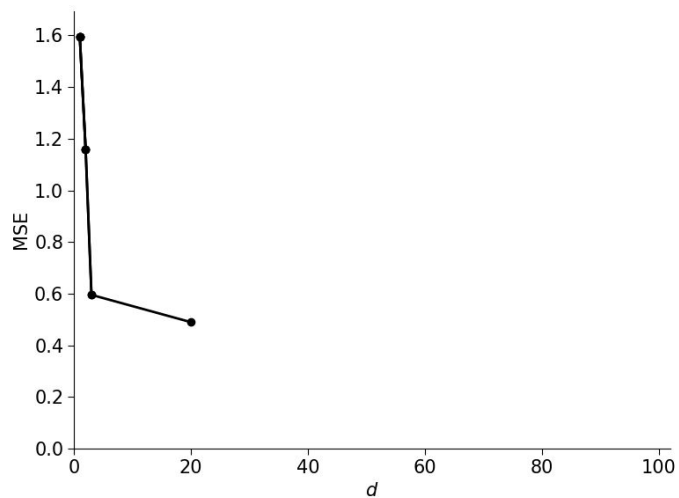
Model complexity



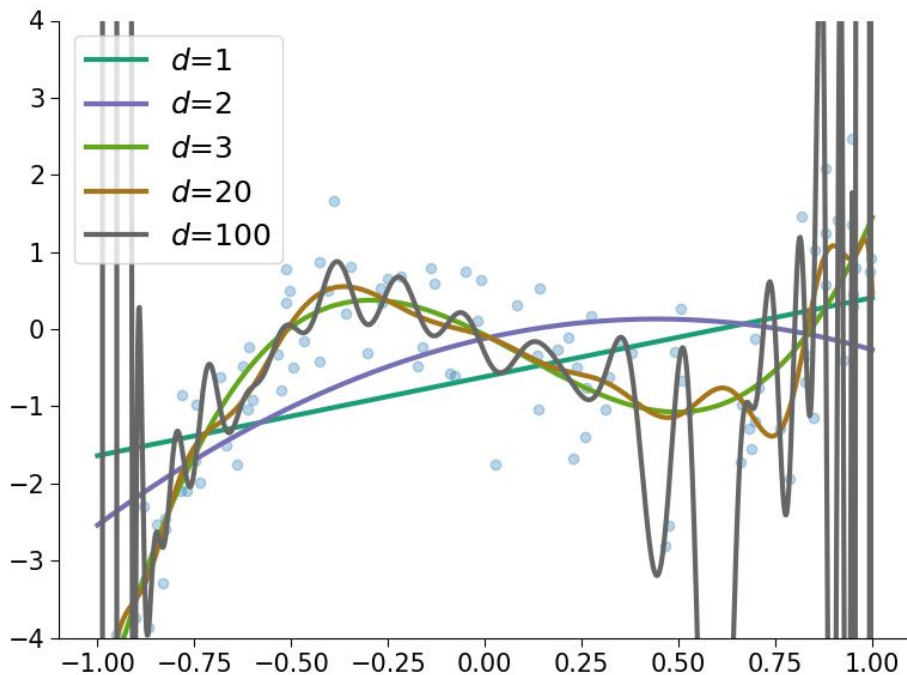
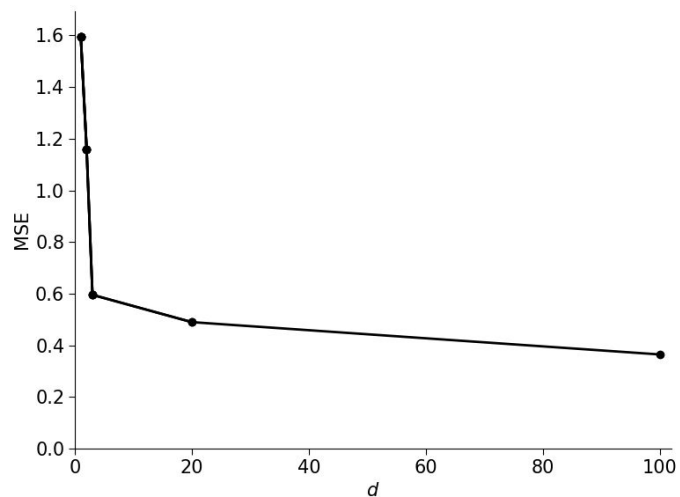
Model complexity



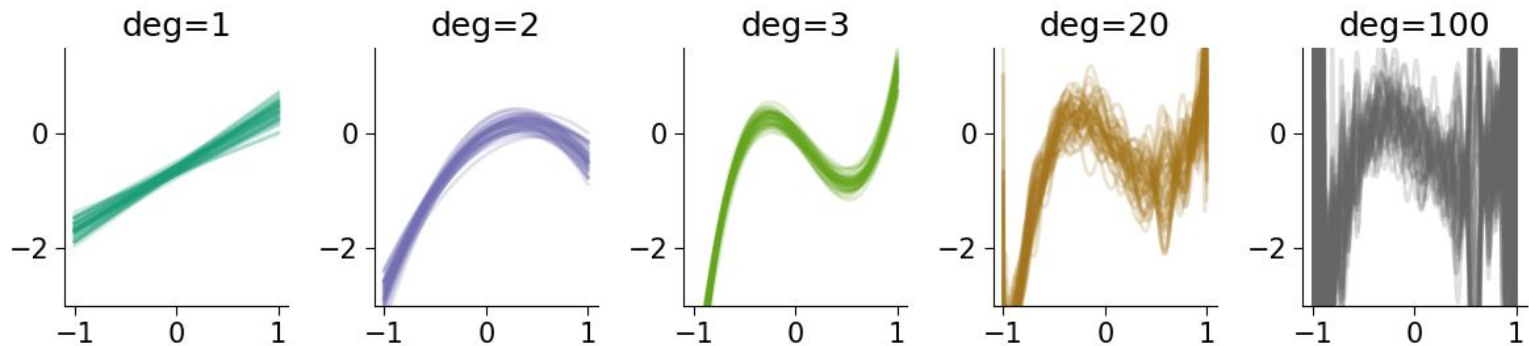
Model complexity



Model complexity



Model complexity



The bias-variance tradeoff

The complex models have:

- Low **bias** – they will fit whatever sample we give them
- High **variance** – different samples will result in very different predictions

The simple models have:

- High **bias** – they will systematically differ from the observations
- Low **Variance** – different samples result in similar predictions

The expected generalization error depends on both bias and variance

Bias-Variance decomposition of MSE

- Assume the true data is generated by $\mathbf{y} = \mathbf{f}(\mathbf{x}) + \boldsymbol{\varepsilon}$
- We estimate $\hat{\mathbf{f}}$ based on a sample $(x^1, y^1, \dots, x^N, y^N)$
- what is the expected error on a **new** datapoint \mathbf{x} ?

$$\mathbb{E}[(\mathbf{y} - \hat{\mathbf{f}})^2] = \mathbb{E}[(\mathbf{f} + \boldsymbol{\varepsilon} - \hat{\mathbf{f}})^2]$$

$$= \mathbb{E}[\boldsymbol{\varepsilon}^2 + \mathbf{f}^2 - 2\mathbf{f}\hat{\mathbf{f}} + \hat{\mathbf{f}}^2] \quad (\boldsymbol{\varepsilon} \text{ uncorrelated with } \hat{\mathbf{f}})$$

$$= \sigma^2 + \mathbf{f}^2 - 2\mathbf{f}\mathbb{E}[\hat{\mathbf{f}}] + \text{Var}[\hat{\mathbf{f}}] + \mathbb{E}^2[\hat{\mathbf{f}}] \quad (\mathbb{E}[X^2] = \text{Var}[X] + \mathbb{E}^2[X])$$

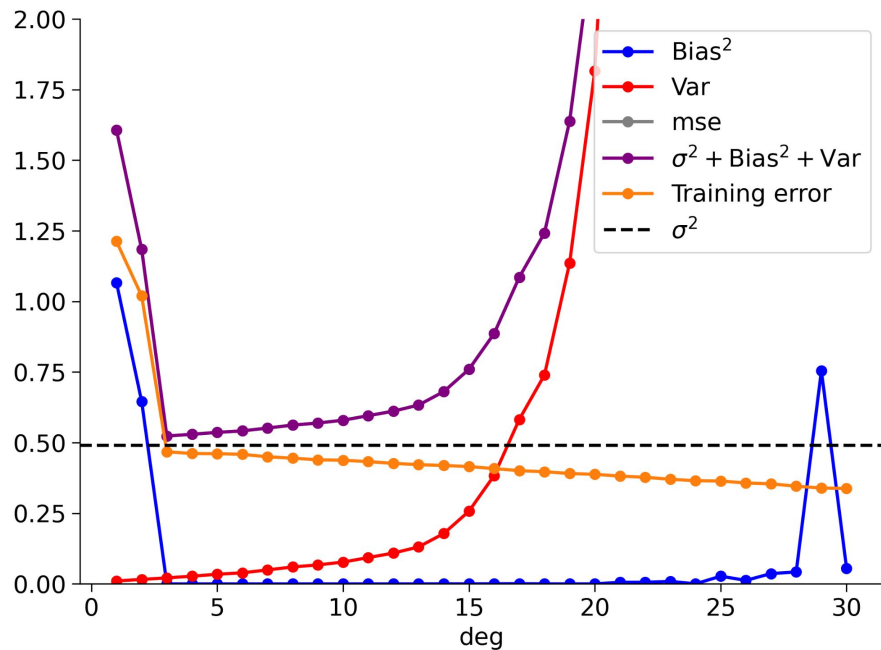
$$= \sigma^2 + (\mathbb{E}[\hat{\mathbf{f}}] - \mathbf{f})^2 + \text{Var}[\hat{\mathbf{f}}]$$

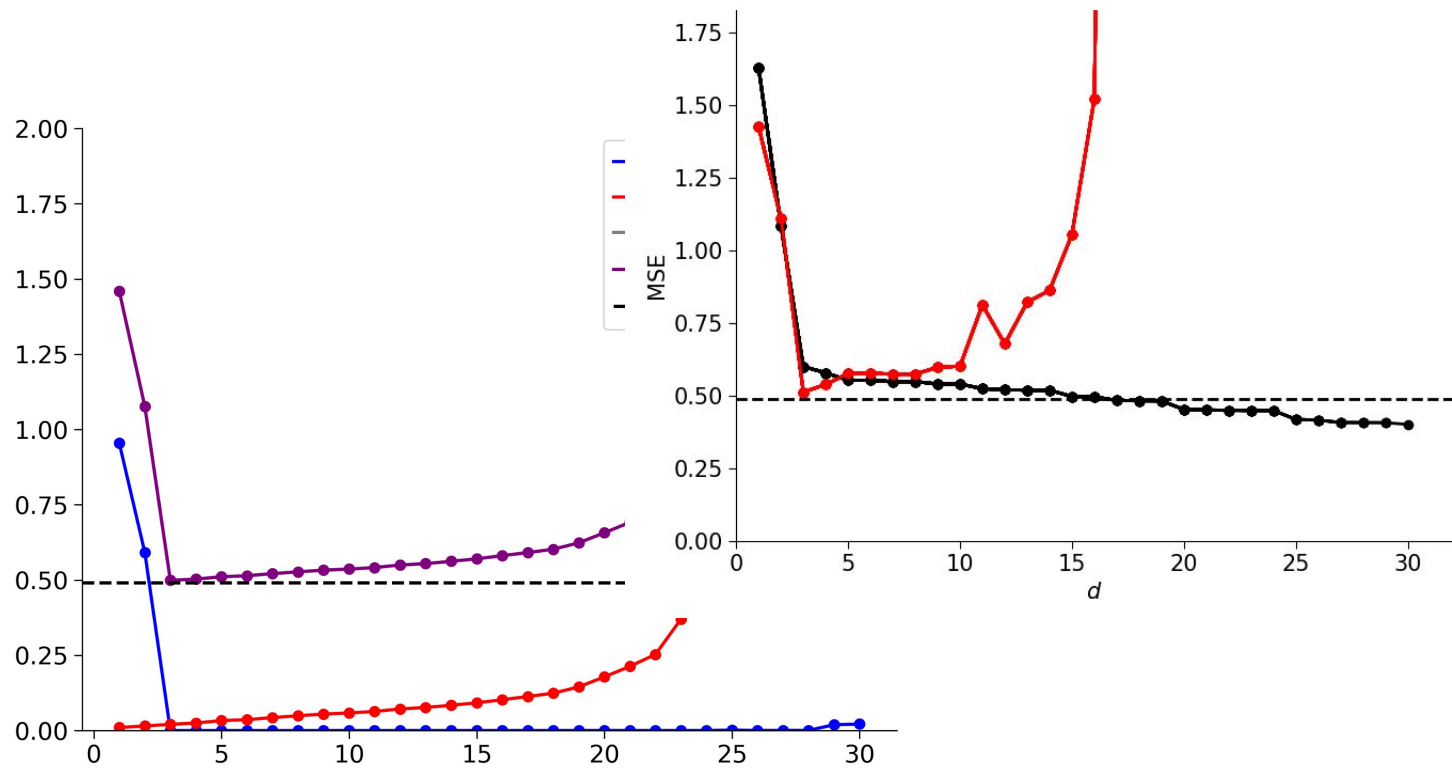
The total expected error: Irreducible error + Bias squared + Variance

The polynomial regression example, again

- Training error keep decreasing
- Generalization error diverges

Making the model more and more complex was **not** the right thing to do





Regularization

- It is sometimes beneficial to explicitly constraint model complexity, even if we make it biased as a result. This is known as regularization
- The typical way: constraint the magnitude of β , “shrinking” parameters towards 0
- Solve a different optimisation problem. For example:

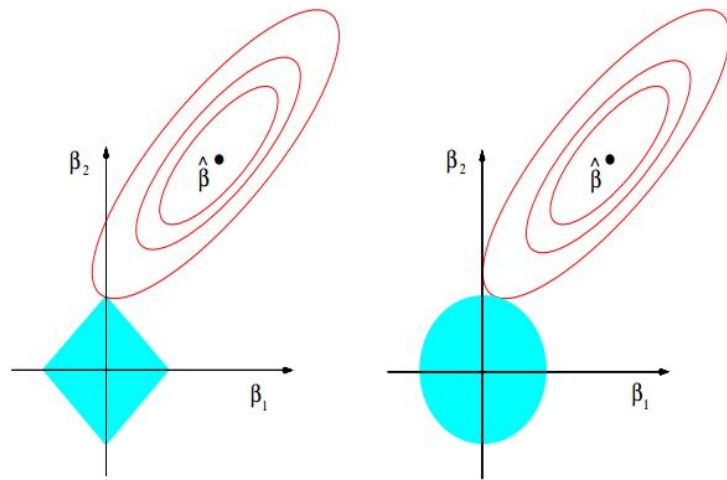
$$\beta = \min\{(y - X\beta)^T(y - X\beta) + \lambda\beta^T\beta\} \quad (\text{Ridge regression})$$

$$\beta = (X^T X + \lambda I)^{-1} X^T y \quad (\text{Solution})$$

- Alternatively, could use L_1 (rather than L_2) norm, resulting in the Lasso objective

Regularization – Ridge and Lasso

- Both methods shrink the parameters, but there are some differences
- Ridge will shrink magnitude of all parameters without setting them to 0
- Lasso will result in a sparse solution, setting small parameters to 0



Ridge regression – geometric intuition

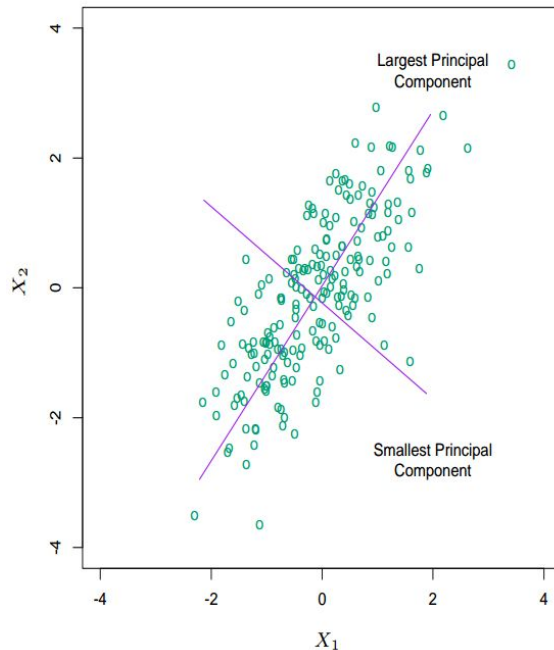
- Recall – the SVD decomposition $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$
- Plugging into the (ordinary) solution:

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

- The orthogonal projection onto the column space of \mathbf{X} , as we've seen before
- The ridge regression solution is

$$\mathbf{X}\mathbf{b} = \mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y} = \mathbf{U}\mathbf{\Sigma}(\mathbf{\Sigma}^2 + \lambda\mathbf{I})^{-1}\mathbf{\Sigma}\mathbf{U}^T\mathbf{y}$$

- Same projection, but each component is scaled down by $d_i^2/(d_i^2 + \lambda)$



Probabilistic interpretation of linear regression

- Back to the basics: why minimise the **square error**?
- If we assume that the data was generated by an unknown linear model with additive gaussian noise, this is equivalent to **maximising the likelihood**:

$$\begin{aligned}\log P(y_1, \dots, y_N | b) &= \log P(y_1 | b) + \dots + \log P(y_N | b) \\ &= -0.5 \sigma^{-2} (y_1 - x_1 b)^2 - \dots - 0.5 \sigma^{-2} (y_N - x_N b)^2 + C\end{aligned}$$

- Therefore, the **maximum-likelihood estimator** for β is the least square estimator

Adding a prior

- We use the data to update our belief about the true parameter.
- But what if we already have some prior belief? We use Bayes' theorem:

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

- For example, we might assume a-priori that $\beta \sim \mathbf{N}(\mathbf{0}, \tau \mathbf{I})$. Then, the log posterior is:

$$\log p(\beta | D) = C - 0.5 \tau^{-1} \beta^T \beta - 0.5 \sigma^{-2} (y - X\beta)^T (y - X\beta)$$

- Maximizing this is the same as solving ridge regression, with $\lambda = \sigma^2 / \tau$
- Consider the extreme cases:

“Uniform”/flat prior: Large τ , resulting in $\lambda \rightarrow 0$, reducing to OLS

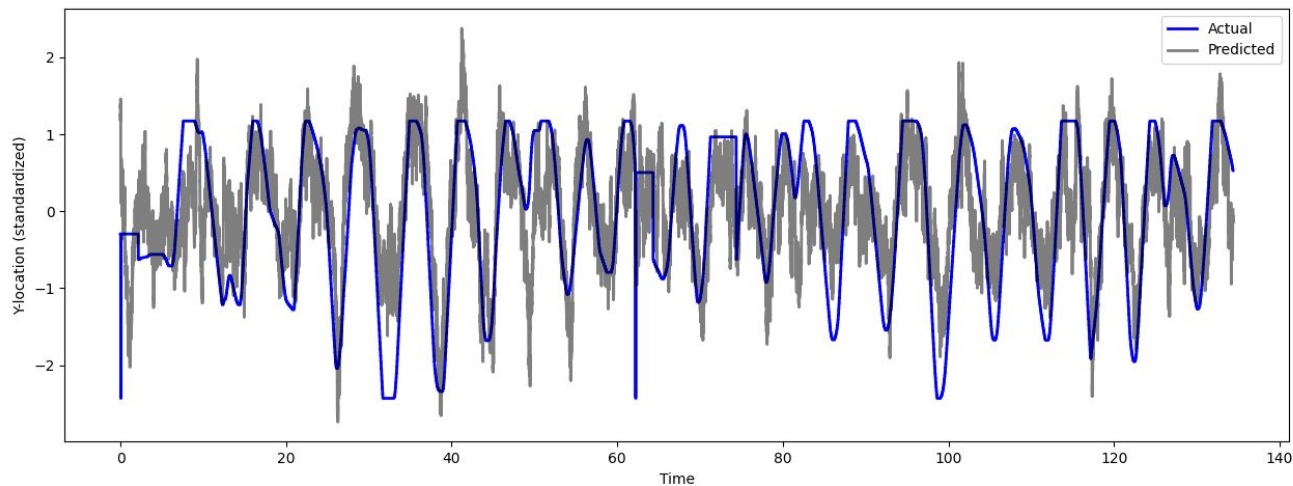
Strong prior: Small τ , resulting in $\lambda \rightarrow \infty$, retaining $\beta=0$

Towards Bayesian Linear Regression

- The Gaussian case is special because of conjugacy
 - posterior remains a gaussian with updated parameters
- In this case, the MAP is also the posterior mean
- The general case (non-zero prior mean, general prior covariance) is also tractable
- Serves as the basic intuition/building blocks for many more Bayesian models

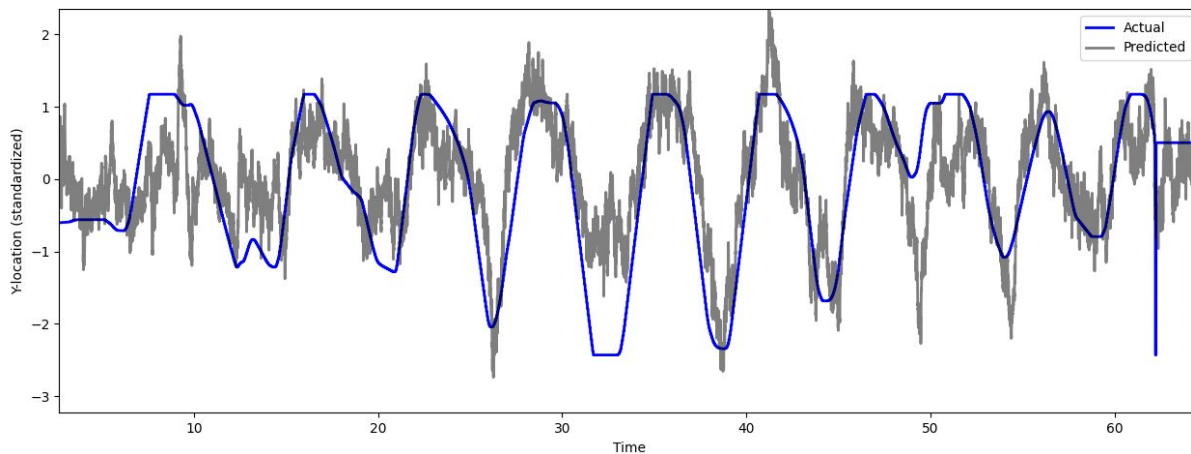
Example 1: Decoding

- We try to predict stimulus/behavior from neural data
- Example: subjects using joystick to move a cursor, following target location
- Predict the cursor location from ECoG data



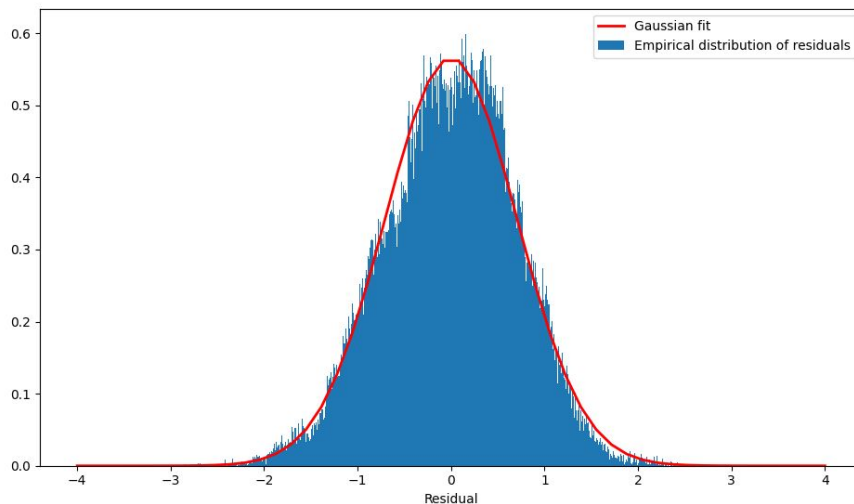
Example 1: Decoding

- We try to predict stimulus/behavior from neural data
- Example: subjects using joystick to move a cursor, following target location
- Predict the cursor location from ECoG data



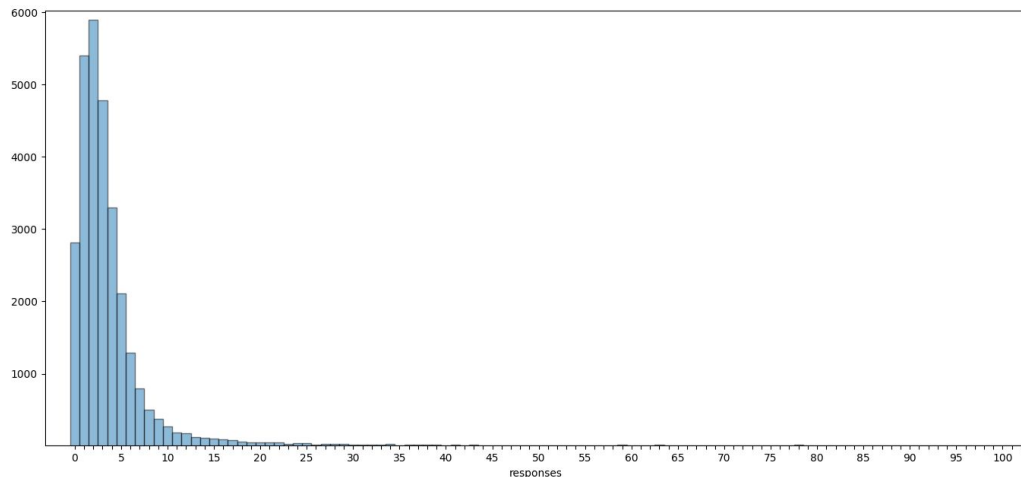
Example 1: Decoding

- We try to predict stimulus/behavior from neural data
- Example: subjects using joystick to move a cursor, following target location
- Predict the cursor location from ECoG data



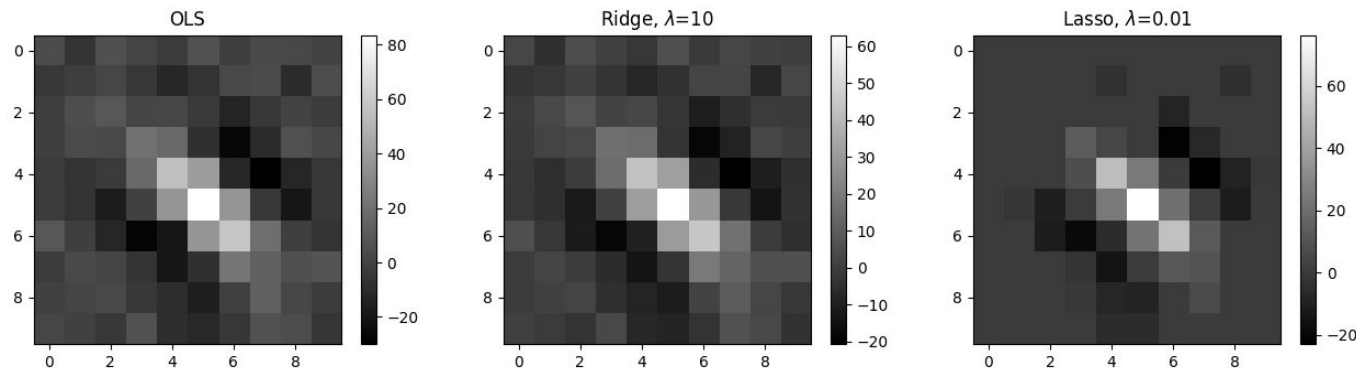
Example 2: Encoding (Receptive Fields)

- Instead of neural→stimulus, we try to predict stimulus→neural
- Example: presented images, simple cell responses (spike counts)
- we really **shouldn't** treat this as ordinary linear regression, but we will anyway



Example 2: Encoding (Receptive Fields)

- Coefficients has the same dimensions as the input – in our case, it's an image
- A way of describing the fitted (linear) receptive field of the cell!



Generalizations to different types of data

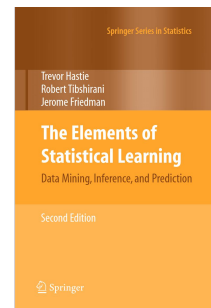
- The assumptions beyond the noise model:
 - Errors are symmetric around 0
 - Small errors are more likely than large errors
- Makes sense for continuous observations corrupted by measurement noise
- What if we have different type of observations?
 - Integers (e.g., number of spikes)
 - Binary outcomes (e.g., behavioral decision)
 - Categorical outcomes
- Different generalization exist (logistic regression, poisson regression, etc)
- Typically, cannot be solved in closed form (but can be optimised efficiently)

Conclusions and further directions

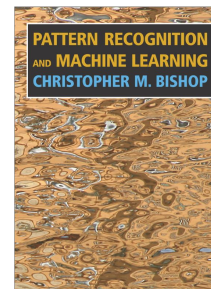
- The model is almost always wrong
- Linear models are simple enough to analyze analytically. This makes them:
 - Extremely useful in their own right
 - Important for building intuition and understanding of more complicated models
- Ultimately, we built a model for **correlations** among predictors and outcomes
- Extra caution should be taken in interpreting the results as a “causal” story, particularly if fitted to observational data
- Many topics we haven’t covered
 - Cross-validation, GLMs, (fully) Bayesian linear regression, subset selection, ...

Sources and materials

The Elements of Statistical Learning, Hastie, Tibshirani, Friedman
(available online)



Pattern recognition and machine learning, Bishop



Nueormatch academy tutorials

<https://compneuro.neuromatch.io/tutorials/intro.html>

LOTS of material online – but not everything is equally good