

Linear prediction: features, overfitting, and losses

CS 446 / ECE 449

2022-02-04 23:03:41 -0600 (1fc7e22)

Plan for today

- ▶ Features.
- ▶ Overfitting.
- ▶ Loss construction and multiclass output.

Recall: solution to Old Faithful example.

1. Form pairs $\mathbf{x}_i := \begin{bmatrix} b_i - a_i \\ 1 \end{bmatrix}$, and matrix

$$\mathbf{X} := \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} = \begin{bmatrix} b_1 - a_1 & 1 \\ & \vdots \\ b_n - a_n & 1 \end{bmatrix} \in \mathbb{R}^{n \times 2}.$$

Form labels $\mathbf{y} \in \mathbb{R}^n$, $y_i := a_{i+1} - b_i$.

2. Choose OLS solution $\hat{\mathbf{w}}_{\text{ols}} := \mathbf{X}^+ \mathbf{y}$.

Recall: solution to Old Faithful example.

1. Form pairs $\mathbf{x}_i := \begin{bmatrix} b_i - a_i \\ 1 \end{bmatrix}$, and matrix

$$\mathbf{X} := \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} = \begin{bmatrix} b_1 - a_1 & 1 \\ & \vdots \\ b_n - a_n & 1 \end{bmatrix} \in \mathbb{R}^{n \times 2}.$$

Form labels $\mathbf{y} \in \mathbb{R}^n$, $y_i := a_{i+1} - b_i$.

2. Choose OLS solution $\hat{\mathbf{w}}_{\text{ols}} := \mathbf{X}^+ \mathbf{y}$.

- ▶ We can view \mathbf{x}_i as the output of a **feature mapping** ϕ :

$$\mathbf{x}_i := \phi(a_i, b_i) := \begin{bmatrix} b_i - a_i \\ 1 \end{bmatrix}.$$

- ▶ Feature mappings are an easy way to make simple methods like linear predictors more powerful.
- ▶ Appending 1 is very common:
instead of the linear rule $\mathbf{x} \mapsto 2 \cdot \mathbb{1}[\mathbf{w}^\top \mathbf{x} \geq 0] - 1$,
we have the **affine decision rule** $\mathbf{x} \mapsto 2 \cdot \mathbb{1}[\mathbf{w}^\top \mathbf{x} \geq -b] - 1$.

“pytorch meta-algorithm”.

1. Clean/augment data (lecture 10?).
2. Pick model/architecture (anything from lectures 2-13).
3. Pick a loss function measuring model fit to data.
4. Run a gradient descent variant to fit model to data.
5. Tweak 1-4 until training error is small.
6. Tweak 1-5, possibly reducing model complexity, until testing error is small.

As part of step 1:

choose $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$ and replace \mathbf{x} with $\phi(\mathbf{x})$ everywhere.

“pytorch meta-algorithm”.

1. Clean/augment data (lecture 10?).
2. Pick model/architecture (anything from lectures 2-13).
3. Pick a loss function measuring model fit to data.
4. Run a gradient descent variant to fit model to data.
5. Tweak 1-4 until training error is small.
6. Tweak 1-5, possibly reducing model complexity, until testing error is small.

As part of step 1:

choose $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$ and replace \mathbf{x} with $\phi(\mathbf{x})$ everywhere.

Remark. Sometimes called feature engineering.

It still exists, despite promises of deep learning.

What if data are not vectors: $\mathbf{x} \in \mathcal{X} \neq \mathbb{R}^d$?

Common example: \mathcal{X} is the set of English words. Two approaches:

1. **Bag-of-words:** pick some ordering on \mathcal{X} , and map the i^{th} word to $e_i \in \mathbb{R}^{|\mathcal{X}|}$.
Can encode documents as normalized word counts.
2. **Word2vec:** magic embedding $\mathcal{X} \rightarrow \mathbb{R}^d$ with $d \ll |\mathcal{X}|$ using deep networks.

Both are still heavily used.

Both have many nuances (e.g., “compressing” \mathcal{X} so that have/having/had/etc all have the same mapping).

Example feature mapping of vector data: monomial features

Suppose $\phi(x) = (1, x, x^2)$, and $\mathbf{w} = (a, b, c)$. Then

$$x \mapsto \mathbf{w}^\top \phi(x) = a + bx + cx^2,$$

and our linear learning tools can now learn quadratic polynomials!

Example feature mapping of vector data: monomial features

Suppose $\phi(x) = (1, x, x^2)$, and $\mathbf{w} = (a, b, c)$. Then

$$x \mapsto \mathbf{w}^\top \phi(x) = a + bx + cx^2,$$

and our linear learning tools can now learn quadratic polynomials!

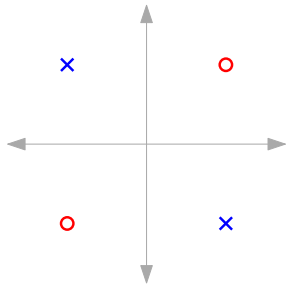
Multivariate case: define

$$\phi(\mathbf{x}) := (1, \underbrace{x_1, \dots, x_d}_{\text{linear terms}}, \underbrace{x_1^2, \dots, x_d^2}_{\text{squared terms}}, \underbrace{x_1x_2, \dots, x_1x_d, \dots, x_{d-1}x_d}_{\text{cross terms}}).$$

Note: hw1 has a coding problem with polynomial expansion, and the order there is slightly different.

XOR example

Consider XOR: $\mathbf{x} \in (\pm 1, \pm 1)$, and $y = x_1x_2$.



Not linearly separable with provided features,
but linearly separable with quadratic features!

Overfitting

“pytorch meta-algorithm”.

⋮

5. Tweak 1-4 until **training error** is small.
6. Tweak 1-5, **possibly reducing model complexity**, until **testing error** is small.

“pytorch meta-algorithm”.

⋮

5. Tweak 1-4 until **training error** is small.
6. Tweak 1-5, **possibly reducing model complexity**, until **testing error** is small.

Feature expansion makes step 5 easy:
given \mathbf{x} , predict using $\phi(\mathbf{x})$, where

$$\phi(\mathbf{x}) = \begin{cases} \mathbf{e}_i & \text{exists } \mathbf{x}_i = \mathbf{x} \text{ in the training data,} \\ 0 & \text{otherwise.} \end{cases}$$

What about step 6?

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

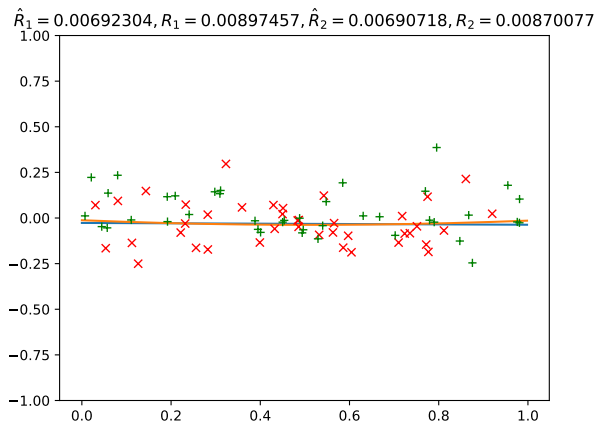
Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



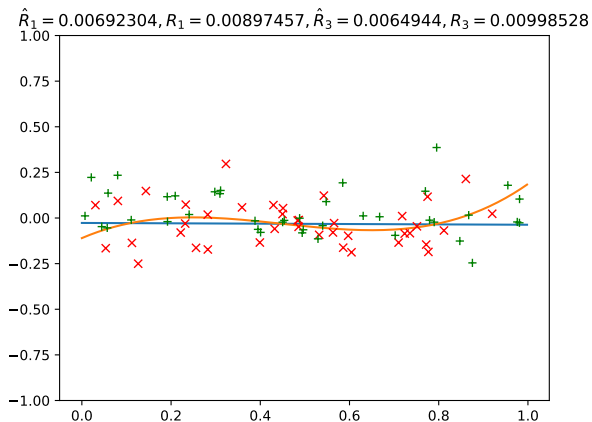
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



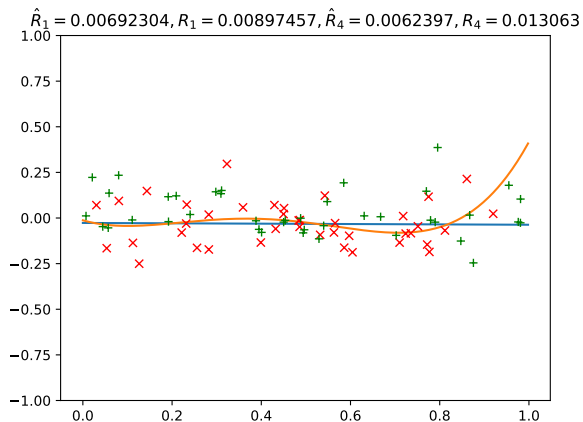
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



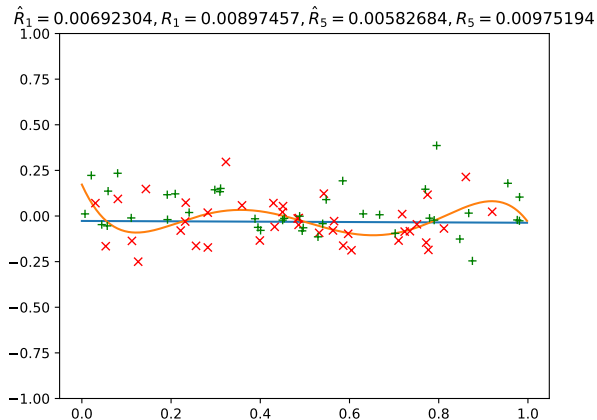
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



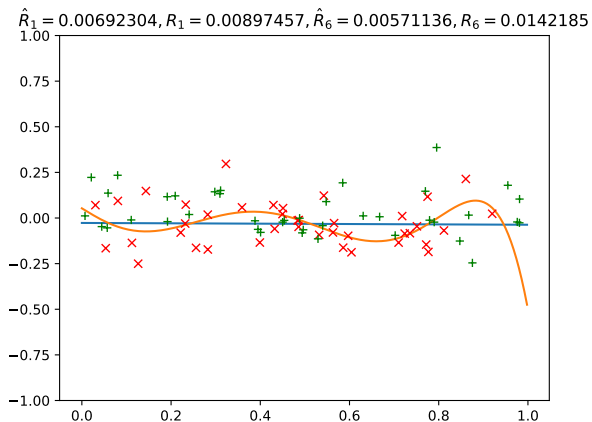
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



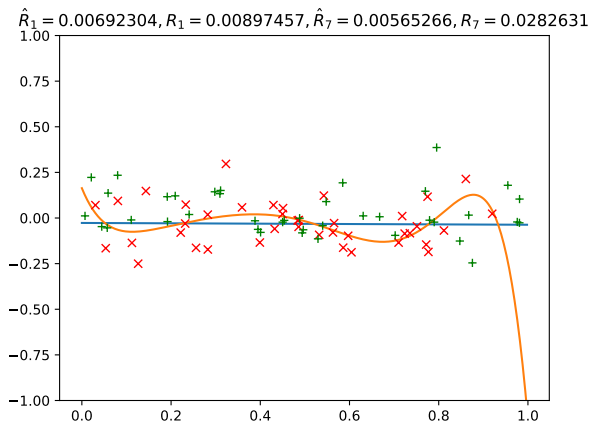
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



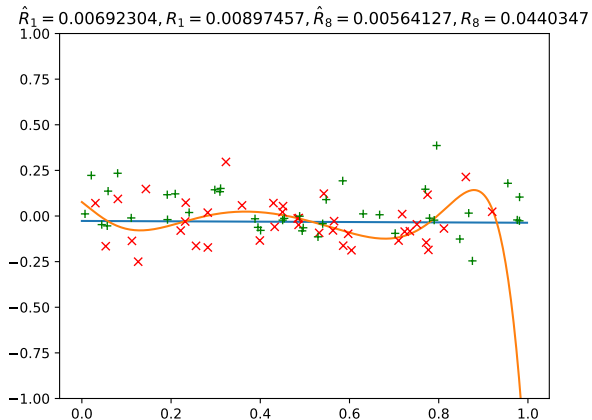
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



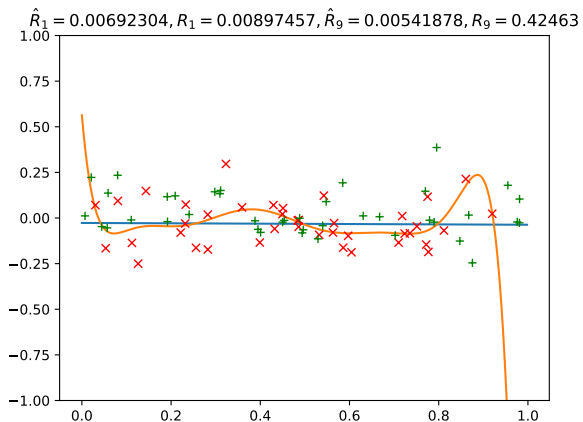
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



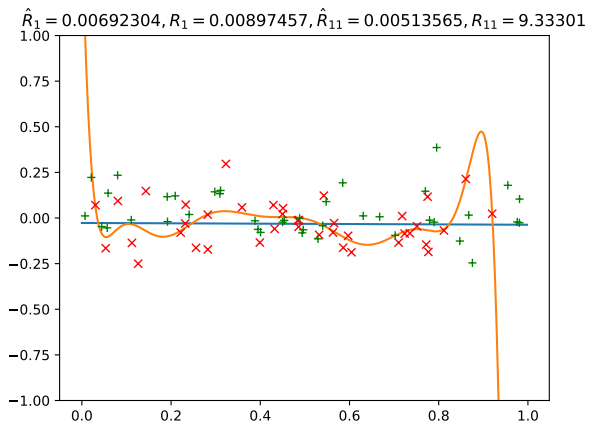
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



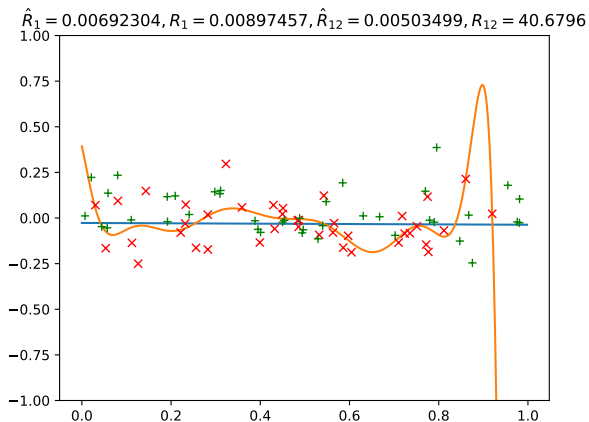
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



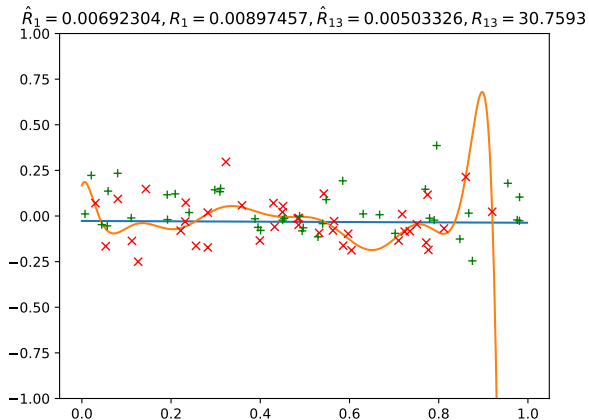
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top(1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



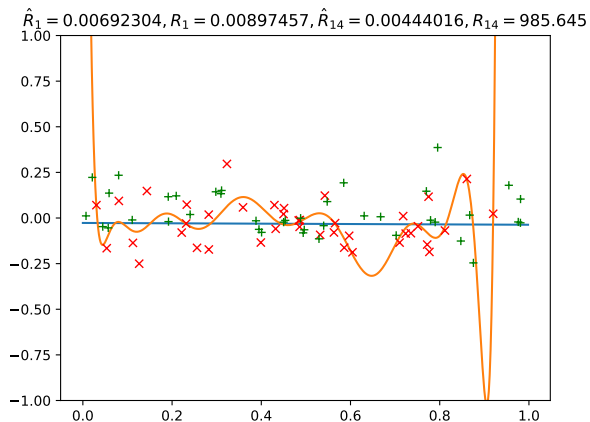
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.

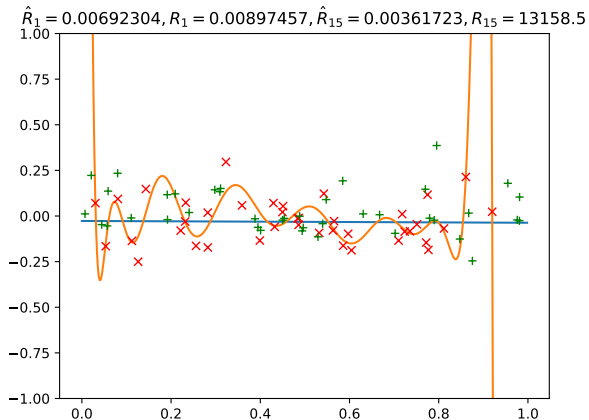


Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



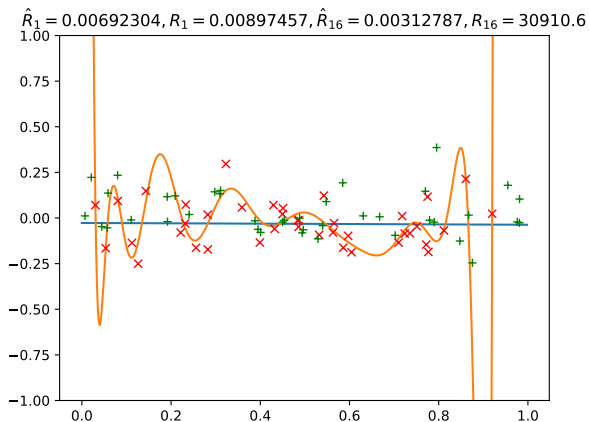
(training data is red, testing data is green.)

Fitting polynomials to a noisy constant function

Data: (x, y) has $x \sim \text{Uniform}([0, 1])$ and $y \sim \mathcal{N}(0, 1)$ (independent!).

Predictor: polynomials of varying degree, $x \mapsto \mathbf{w}^\top (1, x, x^2, \dots, x^r)$.

Method: OLS solution $\hat{\mathbf{w}}_{\text{ols}}$.



(training data is red, testing data is green.)

What happened to our test error?

Degree 1: training error 0.0069, testing error 0.00897.

Degree 16: training error 0.0031, testing error 30910.6.

What happened to our test error?

Degree 1: training error 0.0069, testing error 0.00897.

Degree 16: training error 0.0031, testing error 30910.6.

Overfitting is when training error is good, but testing error is bad.

What happened to our test error?

Degree 1: training error 0.0069, testing error 0.00897.

Degree 16: training error 0.0031, testing error 30910.6.

Overfitting is when training error is good, but testing error is bad.

- ▶ Often it means a model is complicated in a way that is incompatible with the data. Step 6 of pytorch meta-algorithm is typically about tuning model complexity.
- ▶ Sometimes this is called a bias/variance tradeoff: the bias (observed training error) is good, but the variance (random test error) is bad.
- ▶ It is not true that “complex models overfit” (deep networks are often fine).
The situation is not well-understood.

One approach: regularization

Reduce model complexity via **regularized ERM**: pick $\lambda > 0$, and approximately solve

$$\min_{\mathbf{w} \in \mathbb{R}^d} \widehat{\mathcal{R}}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

$\lambda \|\mathbf{w}\|_2^2$ is **regularization**, and this is just one choice.

One approach: regularization

Reduce model complexity via **regularized ERM**: pick $\lambda > 0$, and approximately solve

$$\min_{\mathbf{w} \in \mathbb{R}^d} \widehat{\mathcal{R}}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

$\lambda \|\mathbf{w}\|_2^2$ is **regularization**, and this is just one choice.

- ▶ For least squares, it has a nice form: the solution becomes

$$\left(\mathbf{X}^T \mathbf{X} + \lambda n \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y},$$

where this inverse always exists (why?).

(It has a name: “ridge regression”.)

- ▶ This perspective has lost favor in deep networks, as they can generalize well even if weight norms are large, and regularization doesn't change their test error much.
- ▶ λ is a **hyper-parameter**; needs to be found some other way.
- ▶ Sometimes this specific regularizer is called **weight decay**:

$$\nabla_{\mathbf{w}} \left(\widehat{\mathcal{R}}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right) = \nabla_{\mathbf{w}} \widehat{\mathcal{R}}(\mathbf{w}) + \lambda \mathbf{w}.$$

Loss functions and multiclass prediction

Loss functions and multiclass prediction

So far we have seen:

$$\ell_{\text{sq}}(z) = \frac{1}{2}(1 - z)^2, \quad (\text{squared loss}),$$

$$\ell_{\text{logistic}}(z) = \ln(1 + \exp(-z)), \quad (\text{logistic loss}).$$

Loss functions and multiclass prediction

So far we have seen:

$$\ell_{\text{sq}}(z) = \frac{1}{2}(1 - z)^2, \quad (\text{squared loss}),$$

$$\ell_{\text{logistic}}(z) = \ln(1 + \exp(-z)), \quad (\text{logistic loss}).$$

Questions:

- ▶ Are there other choices?
- ▶ What about other output spaces, e.g., multiclass?

Loss functions and multiclass prediction

So far we have seen:

$$\ell_{\text{sq}}(z) = \frac{1}{2}(1 - z)^2, \quad (\text{squared loss}),$$

$$\ell_{\text{logistic}}(z) = \ln(1 + \exp(-z)), \quad (\text{logistic loss}).$$

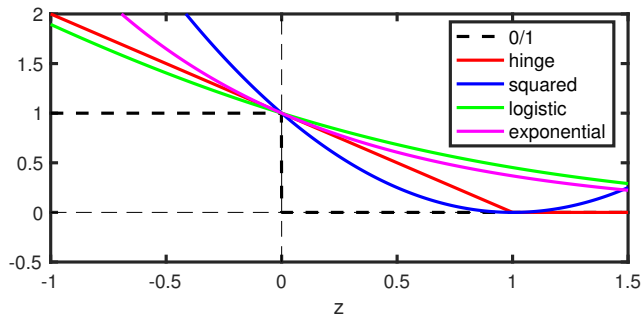
Questions:

- ▶ Are there other choices?
- ▶ What about other output spaces, e.g., multiclass?

Remark: We introduced squared loss as $\ell_{\text{sq}}(\hat{y}, y) = (y - \hat{y})^2/2$.

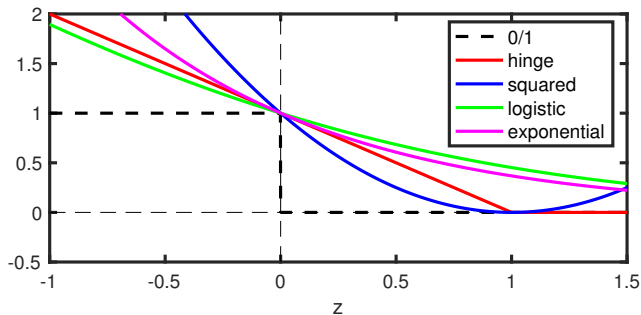
This makes sense for regression; for classification, combining arguments into a single $z := \hat{y}y$ suffices.

Standard classification losses



hinge	$l_{\text{hinge}}(z) = \max\{0, 1 - z\}$
squared	$2l_{\text{sq}}(z) = (1 - z)^2$
logistic	$l_{\text{logistic}}(z) / \ln(2) = \ln(1 + e^{-z}) / \ln(2)$
exponential	$l_{\text{exp}}(z) = e^{-z}$

Standard classification losses



hinge	$\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$
squared	$2\ell_{\text{sq}}(z) = (1 - z)^2$
logistic	$\ell_{\text{logistic}}(z) / \ln(2) = \ln(1 + e^{-z}) / \ln(2)$
exponential	$\ell_{\text{exp}}(z) = e^{-z}$

- ▶ Upper bound zero-one loss: $\ell(z) \rightarrow 0$ implies $\mathbb{1}[z \leq 0] \rightarrow 0$.
Sometimes called **convex surrogate losses** (for the zero-one loss).
- ▶ Differentiable and have $\ell'(0) < 0$: helps GD.

Designing losses with maximum likelihood

Stronger goal than classification:
suppose our parameters w define a **conditional model** $p_w(y|\mathbf{x})$.

Designing losses with maximum likelihood

Stronger goal than classification:

suppose our parameters w define a **conditional model** $p_w(y|\mathbf{x})$.

Maximum likelihood (MLE) approach to model training:

given $((\mathbf{x}_i, y_i))_{i=1}^n$, choose w via

$$\arg \max_w \prod_{i=1}^n p_w(y_i|\mathbf{x}_i) = \arg \min_w -\ln \prod_{i=1}^n p_w(y_i|\mathbf{x}_i) = \arg \min_w \sum_{i=1}^n \ln \frac{1}{p_w(y_i|\mathbf{x}_i)}.$$

Designing losses with maximum likelihood

Stronger goal than classification:

suppose our parameters \mathbf{w} define a **conditional model** $p_{\mathbf{w}}(y|\mathbf{x})$.

Maximum likelihood (MLE) approach to model training:

given $((\mathbf{x}_i, y_i))_{i=1}^n$, choose \mathbf{w} via

$$\arg \max_{\mathbf{w}} \prod_{i=1}^n p_{\mathbf{w}}(y_i|\mathbf{x}_i) = \arg \min_{\mathbf{w}} -\ln \prod_{i=1}^n p_{\mathbf{w}}(y_i|\mathbf{x}_i) = \arg \min_{\mathbf{w}} \sum_{i=1}^n \ln \frac{1}{p_{\mathbf{w}}(y_i|\mathbf{x}_i)}.$$

Looks like ERM with loss $\ln \frac{1}{p_{\mathbf{w}}(y|\mathbf{x})}$!

Note: can use these losses even if true label probabilities have a different model.

Squared loss via MLE

Define conditional model $p_{\mathbf{w}}$ using standard Gaussian with mean $\mathbf{x}^T \mathbf{w}$:

$$p_{\mathbf{w}}(y|\mathbf{x}) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y - \mathbf{x}^T \mathbf{w})^2}{2}\right).$$

Squared loss via MLE

Define conditional model $p_{\mathbf{w}}$ using standard Gaussian with mean $\mathbf{x}^T \mathbf{w}$:

$$p_{\mathbf{w}}(y|\mathbf{x}) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y - \mathbf{x}^T \mathbf{w})^2}{2}\right).$$

Then

$$\ln \frac{1}{p_{\mathbf{w}}(y|\mathbf{x})} = \frac{1}{2} (y - \mathbf{x}^T \mathbf{w})^2 + \frac{1}{2} \ln(2\pi),$$

and

$$\begin{aligned} \arg \min_{\mathbf{w}} \sum_{i=1}^n \ln \frac{1}{p_{\mathbf{w}}(y_i|\mathbf{x}_i)} &= \arg \min_{\mathbf{w}} \sum_{i=1}^n \ell_{\text{sq}}(\mathbf{x}_i^T \mathbf{w}, y_i) \\ &= \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell_{\text{sq}}(\mathbf{x}_i^T \mathbf{w}, y_i). \end{aligned}$$

Logistic loss via MLE

Define conditional model

$$p_{\mathbf{w}}(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x}^T \mathbf{w})}.$$

whereby

$$p_{\mathbf{w}}(y = -1|\mathbf{x}) = 1 - p_{\mathbf{w}}(y = 1|\mathbf{x}) = \frac{\exp(-\mathbf{x}^T \mathbf{w})}{1 + \exp(-\mathbf{x}^T \mathbf{w})} = \frac{1}{1 + \exp(\mathbf{x}^T \mathbf{w})}.$$

Logistic loss via MLE

Define conditional model

$$p_{\mathbf{w}}(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x}^T \mathbf{w})}.$$

whereby

$$p_{\mathbf{w}}(y = -1|\mathbf{x}) = 1 - p_{\mathbf{w}}(y = 1|\mathbf{x}) = \frac{\exp(-\mathbf{x}^T \mathbf{w})}{1 + \exp(-\mathbf{x}^T \mathbf{w})} = \frac{1}{1 + \exp(\mathbf{x}^T \mathbf{w})}.$$

Then

$$\begin{aligned} \ln \frac{1}{p_{\mathbf{w}}(y|\mathbf{x})} &= \ln \frac{1}{p_{\mathbf{w}}(y = 1|\mathbf{x})^{(1+y)/2} (1 - p_{\mathbf{w}}(y = 1|\mathbf{x}))^{(1-y)/2}} \\ &= \frac{1+y}{2} \ln(1 + \exp(-\mathbf{x}^T \mathbf{w})) + \frac{(1-y)}{2} \ln(1 + \exp(\mathbf{x}^T \mathbf{w})) \\ &= \ln(1 + \exp(-y\mathbf{x}^T \mathbf{w})), \end{aligned}$$

and

$$\arg \min_{\mathbf{w}} \sum_{i=1}^n \ln \frac{1}{p_{\mathbf{w}}(y_i|\mathbf{x}_i)} = \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell_{\text{logistic}}(y_i \mathbf{x}_i^T \mathbf{w}).$$

Multiclass classification

Setup: given \mathbf{x} predict $y \in \{1, \dots, k\}$.

Two approaches:

- ▶ **One-vs-all** (a reduction approach): train k binary classifiers, each with $\tilde{y}_i := 2 \cdot \mathbb{1}[y_i = j] - 1$. Predict by choosing the largest output.
- ▶ Directly train a model $p(y = j | \mathbf{x})$ with maximum likelihood.

Cross-entropy loss (multi-class logistic loss) via MLE

Conditional model: given predictor $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$, model $p_f(\hat{y} = \cdot | \mathbf{x}) \propto \exp(f(\mathbf{x}))$,
which means $p_f(\hat{y} = j | \mathbf{x}) = \frac{\exp(f(\mathbf{x})_j)}{\sum_{i=1}^k \exp(f(\mathbf{x})_i)}$.

Corresponding **cross-entropy loss** ℓ_{ce} : given true label $y \in \{1, \dots, k\}$,

Cross-entropy loss (multi-class logistic loss) via MLE

Conditional model: given predictor $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$, model $p_f(\hat{y} = \cdot | \mathbf{x}) \propto \exp(f(\mathbf{x}))$,
which means $p_f(\hat{y} = j | \mathbf{x}) = \frac{\exp(f(\mathbf{x})_j)}{\sum_{i=1}^k \exp(f(\mathbf{x})_i)}$.

Corresponding **cross-entropy loss** ℓ_{ce} : given true label $y \in \{1, \dots, k\}$,

$$\begin{aligned} \ln \frac{1}{p_f(\hat{y} = y | \mathbf{x})} &= \ln \frac{\sum_{j=1}^k \exp(f(\mathbf{x})_j)}{\exp(f(\mathbf{x}))_y} \\ &= -f(\mathbf{x})_y + \ln \sum_{j=1}^k \exp(f(\mathbf{x})_j) \\ &=: \ell_{\text{ce}}(f(\mathbf{x}), y). \end{aligned}$$

Cross-entropy loss (multi-class logistic loss) via MLE

Conditional model: given predictor $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$, model $p_f(\hat{y} = \cdot | \mathbf{x}) \propto \exp(f(\mathbf{x}))$,
which means $p_f(\hat{y} = j | \mathbf{x}) = \frac{\exp(f(\mathbf{x})_j)}{\sum_{i=1}^k \exp(f(\mathbf{x})_i)}$.

Corresponding **cross-entropy loss** ℓ_{ce} : given true label $y \in \{1, \dots, k\}$,

$$\begin{aligned} \ln \frac{1}{p_f(\hat{y} = y | \mathbf{x})} &= \ln \frac{\sum_{j=1}^k \exp(f(\mathbf{x})_j)}{\exp(f(\mathbf{x})_y)} \\ &= -f(\mathbf{x})_y + \ln \sum_{j=1}^k \exp(f(\mathbf{x})_j) \\ &=: \ell_{ce}(f(\mathbf{x}), y). \end{aligned}$$

Notes.

In pytorch, this is `torch.nn.CrossEntropyLoss()` ($f(\mathbf{x})$, y).

Loss is minimized when $p_f(\hat{y} = y | \mathbf{x}) \approx 1$.

Name comes from the cross-entropy expression $\sum_{j=1}^k (e_y)_j \ln \frac{1}{p_f(\hat{y}=j|\mathbf{x})}$.

Summary for today

- ▶ Features.
- ▶ Overfitting.
- ▶ Loss construction and multiclass output.

(Appendix.)

FPR/TPR, confusion matrix, ROC/AUC curve, cost-sensitive losses.
KL and cross-entropy.

Multivariate prediction: regression case

Suppose we have k labels for each example:

$$\mathbf{Y} := \begin{bmatrix} \leftarrow & \mathbf{y}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{y}_n^\top & \rightarrow \end{bmatrix}.$$

It's natural to also learn $\mathbf{W} \in \mathbb{R}^{d \times k}$:

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{d \times k}} \|\mathbf{XW} - \mathbf{Y}\|_F^2.$$

Multivariate prediction: regression case

Suppose we have k labels for each example:

$$\mathbf{Y} := \begin{bmatrix} \leftarrow & \mathbf{y}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{y}_n^\top & \rightarrow \end{bmatrix}.$$

It's natural to also learn $\mathbf{W} \in \mathbb{R}^{d \times k}$:

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{d \times k}} \|\mathbf{XW} - \mathbf{Y}\|_F^2.$$

This is the same as k regular linear regressions: given \mathbf{W} ,

$$\|\mathbf{XW} - \mathbf{Y}\|_F^2 = \sum_{j=1}^n \|\mathbf{XW}_{:j} - \mathbf{Y}_{:j}\|^2 = \sum_{j=1}^n \sum_{i=1}^n \left(\mathbf{X}_{i:}^\top \mathbf{W}_{:j} - y_{ij} \right)^2.$$

We have **reduced** multivariate regression to univariate regression:

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times k}} \|\mathbf{XW} - \mathbf{Y}\|_F^2 = \sum_{j=1}^k \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{Xw} - \mathbf{Y}_{:j}\|^2.$$

Remark on modeling assumptions (incomplete slide!)

1. It is **not true** that logistic loss and squared loss only work well when the preceding likelihood models are correct.
2. There is a complicated way to argue that they can still behave optimally even when their model is wrong, but here I'll give a simpler perspective.
3. We can **derive** the notion of sample mean and sample covariance by solving for mean and variance in a maximum likelihood computation for multivariate Gaussian parameters over some point set. That is to say, the sample mean and covariance will recover the Gaussian parameters **when the data is generated by a Gaussian**. However, **when the data is not generated by a Gaussian**, while we are not "recovering parameters of an underlying model", we are still obtaining useful statistics which describe the data well (and, e.g., if we were to get all moments, we would in a sense completely characterize the data).

- ▶ Shalev-Shwartz/Ben-David: chapters 9, 24.