Supervised Learning via Empirical Risk Minimization

Jong-Han Kim

EE787 Fundamentals of machine learning Kyung Hee University

Predictors

Data fitting

lacksimwe think $y\in {f R}$ and $x\in {f R}^d$ are (approximately) related by

 $y \approx f(x)$

- ▶ *x* is called the *independent variable* or *feature vector*
- ▶ y is called the outcome or response or target or label or dependent variable
- \blacktriangleright often y is something we want to predict
- \blacktriangleright we don't know the 'true' relationship between x and y

Features

often x is a vector of features:

- ▶ documents
 - ▶ x is word count histogram for a document
- ▶ patient data
 - \blacktriangleright x are patient attributes, test results, symptoms
- customers
 - $\triangleright x$ is purchase history and other attributes of a customer

Where features come from

- we use u to denote the raw input data, such as a vector, word or text, image, video, audio, ...
- $x = \phi(u)$ is the corresponding *feature vector*
- **•** the function ϕ is called the *embedding* or *feature function*
- $\blacktriangleright \phi$ might be very simple or quite complicated
- \blacktriangleright similarly, the raw output data v can be featurized as $y=\psi(v)$
- \blacktriangleright often we take $\phi(u)_1 = x_1 = 1$, the *constant feature*
- (much more on these ideas later)

Data and prior knowledge

- \blacktriangleright we are given data $x^1,\ldots,x^n\in\mathsf{R}^d$ and $y^1,\ldots,y^n\in\mathsf{R}$
- \blacktriangleright (x^i, y^i) is the *i*th *data pair* or *observation* or *example*
- ▶ we also (might) have *prior knowledge* about what f might look like
 - \blacktriangleright e.g., f is smooth or continuous: $f(x) pprox f(ilde{x})$ when x is near $ilde{x}$
 - or we might know $y \ge 0$

Predictor

- we seek a *predictor* or *model* $g : \mathbb{R}^d \to \mathbb{R}$
- ▶ for feature vector x, our prediction (of y) is $\hat{y} = g(x)$
- predictor g is chosen based on both data and prior knowledge
- ▶ in terms of raw data, our predictor is

$$\hat{v}=\psi^{-1}(g(\phi(u)))$$

(with a slight variation when ψ is not invertible)

- $igstarrow \hat{y}^i pprox y^i$ means our predictor does well on ith data pair
- **>** but our real goal is to have $\hat{y} \approx y$ for (x, y) pairs we have not seen

Information flow



Prediction methods

- fraud, psychic powers, telepathy, magic sticks, incantations, crystals, hunches, statistics, AI, machine learning, data science
- and many algorithms . . .
- example: nearest neighbor predictor
 - **•** given x, find its nearest neighbor x^i among given data
 - \blacktriangleright then predict $\hat{y} = g(x) = y^i$

A learning algorithm is a recipe for producing a predictor given data

Example: Nearest neighbor prediction



- left plot shows nearest neighbor prediction
- ▶ right plot shows fit with cubic polynomial

Linear predictors

Linear predictor

 \blacktriangleright predictors that are linear functions of x are widely used

a linear predictor has the form

$$g(x) = \theta^{\mathsf{T}} x$$

for some vector $\theta \in \mathbf{R}^d$, called the *predictor parameter vector*

▶ also called a *regression model*

 \triangleright x_j is the *j*th feature, so the prediction is a linear combination of features

$$\hat{y}=g(x)= heta_1x_1+\dots+ heta_dx_d$$

 \blacktriangleright we get to choose the predictor parameter vector $heta \in \mathbf{R}^d$

> sometimes we write $g_{\theta}(x)$ to emphasize the dependence on θ

Interpreting a linear predictor

$$\hat{y}=g(x)= heta_1x_1+\dots+ heta_dx_d$$

- θ₃ is the amount that prediction ŷ = g(x) increases when x₃ increases by 1

 particularly interpretable when x₃ is Boolean (only takes values 0 or 1)

 θ₇ = 0 means that the prediction does not depend on x₇
- \blacktriangleright θ small means predictor is insensitive to changes in x:

$$|g(x)-g(ilde{x})|=\left| heta^{ op}x- heta^{ op} ilde{x}
ight|=\left| heta^{ op}(x- ilde{x})
ight|\leq || heta||\;||x- ilde{x}||$$

Affine predictor

- \blacktriangleright suppose the first feature is constant, $x_1=1$
- ▶ the linear predictor g is then an *affine function* of $x_{2:d}$, *i.e.*, linear plus a constant

$$g(x)= heta^{ op}x= heta_1+ heta_2x_2+\dots+ heta_dx_d$$

- θ_1 is called the *offset* or *constant term* in the predictor
- \triangleright θ_1 is the prediction when all features (except the constant) are zero

Empirical risk minimization

Loss function

a loss or risk function $\ell: R \times R \to R$ quantifies how well (more accurately, how badly) \hat{y} approximates y

- \blacktriangleright smaller values of $\ell(\hat{y},y)$ indicate that \hat{y} is a good approximation of y
- typically $\ell(y,y) = 0$ and $\ell(\hat{y},y) \ge 0$ for all \hat{y}, y

examples

- quadratic loss: $\ell(\hat{y}, y) = (\hat{y} y)^2$
- ▶ absolute loss: $\ell(\hat{y}, y) = |\hat{y} y|$

Empirical risk

how well does the predictor g fit a data set $(x^i, y^i), \ i=1,\ldots,n$, with loss ℓ ?

▶ the *empirical risk* is the average loss over the data points,

$$\mathcal{L}=rac{1}{n}\sum_{i=1}^n\ell(\hat{y}^i,y^i)=rac{1}{n}\sum_{i=1}^n\ell(g(x^i),y^i)$$

 \blacktriangleright if ${\cal L}$ is small, the predictor predicts the given data well

 \blacktriangleright when the predictor is parametrized by θ , we write

$$\mathcal{L}(heta) = rac{1}{n}\sum_{i=1}^n \ell(g_ heta(x^i),y^i)$$

to show the dependence on the predictor parameter θ

Mean square error

▶ for square loss $\ell(\hat{y}, y) = (\hat{y} - y)^2$, empirical risk is *mean-square error* (MSE)

$$\mathcal{L} = \mathsf{MSE} = rac{1}{n}\sum_{i=1}^n (g(x^i) - y^i)^2$$

▶ often we use root-mean-square error, RMSE = $\sqrt{\text{MSE}}$, which has same units/scale as outcomes y^i

Mean absolute error

 \blacktriangleright for absolute value $\ell(\hat{y},y) = |\hat{y} - y|$, empirical risk is *mean-absolute error*

$$\mathcal{L} = rac{1}{n}\sum_{i=1}^n |g(x^i)-y^i|$$

- \blacktriangleright has same units/scale as outcomes y^i
- similar to, but not the same as, mean-square error

Empirical risk minimization

- choosing the parameter θ in a parametrized predictor $g_{\theta}(x)$ is called *fitting* the predictor (to data)
- empirical risk minimization (ERM) is a general method for fitting a parametrized predictor
- **ERM**: choose θ to minimize empirical risk $\mathcal{L}(\theta)$
- \blacktriangleright thus, ERM chooses θ by attempting to match given data
- often there is no analytic solution to this minimization problem, so we use *numerical optimization* to find θ that minimizes $\mathcal{L}(\theta)$ (more on this topic later)