



Department of Computer Science
UNIVERSITY OF COLORADO BOULDER



Machine Learning: Chenhao Tan

University of Colorado Boulder
LECTURE 7

Slides adapted from Noah Smith, Chris Ketelsen

Logistics

- HW1 available on Github, due today

Learning objectives

- Understand why features matter
- Understand feature engineering techniques

Outline

Features matter

Feature engineering techniques

Outline

Features matter

Feature engineering techniques

scary words

Outline of CSCI 4622

We've already covered stuff in **blue**!

- Problem formulations: **classification**, regression
- Supervised techniques: **decision trees**, **nearest neighbors**, **perceptron**, linear models, neural networks, support vector machine, kernel methods
- Unsupervised techniques: clustering, linear dimensionality reduction
- “Meta-techniques”: ensembles, expectation-maximization
- Understanding ML: **limits of learning**, practical issues, bias & fairness
- Recurring themes: (stochastic) gradient descent

Today: (More) Best Practices

You already know:

- Separating training and test data
- Hyperparameter tuning on development data

Understanding machine learning is partly about knowing algorithms and partly about the art of mapping abstract problems to learning tasks.

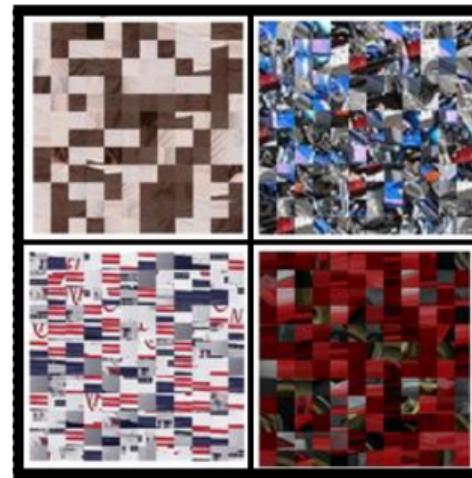
Features Matter

Pixel representation



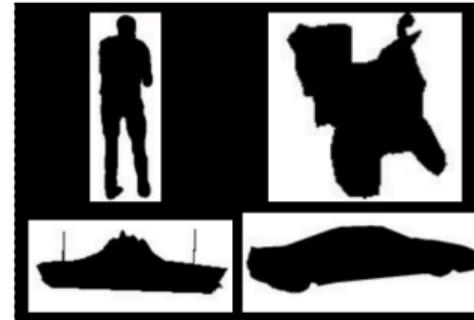
Features Matter

Patch representation



Features Matter

Shape representation



Features Matter

Original



Features Matter

Bag of words

a, algorithms, and, applications, arthur, artificial, building, by, can, coined, computational, computer, computing, construction, data, decisions, designing, detection, difficult, driven, email, employed, evolved, example, explicit, explores, filtering, following, from, good, in, include, infeasible, inputs, instructions, intelligence, intruders, is, learn, learning, machine, make, making, model, name, network, of, on, or, overcome, pattern, performance, predictions, program, programming, range, recognition, sample, samuel, static, strictly, study, such, tasks, that, the, theory, through, vision, was, where, with

Features Matter

Original

The name machine learning was coined in 1959 by Arthur Samuel.[1] Evolved from the study of pattern recognition and computational learning theory in artificial intelligence,[3] machine learning explores the study and construction of algorithms that can learn from and make predictions on data[4] - such algorithms overcome following strictly static program instructions by making data-driven predictions or decisions,[5]:2 through building a model from sample inputs. Machine learning is employed in a range of computing tasks where designing and programming explicit algorithms with good performance is difficult or infeasible; example applications include email filtering, detection of network intruders, and computer vision.

Outline

Features matter

Feature engineering techniques

Irrelevant Features

$$\mathbb{E}[f \mid Y] = \mathbb{E}[f]$$

Irrelevant Features

$$\mathbb{E}[f \mid Y] = \mathbb{E}[f]$$

One irrelevant feature isn't a big deal; what we're worried about is when irrelevant features *outnumber* useful ones!

Irrelevant Features

$$\mathbb{E}[f \mid Y] = \mathbb{E}[f]$$

One irrelevant feature isn't a big deal; what we're worried about is when irrelevant features *outnumber* useful ones!

- Decision trees (not too deep)?

Irrelevant Features

$$\mathbb{E}[f \mid Y] = \mathbb{E}[f]$$

One irrelevant feature isn't a big deal; what we're worried about is when irrelevant features *outnumber* useful ones!

- Decision trees (not too deep)?
Somewhat protected, but beware spurious correlations!

Irrelevant Features

$$\mathbb{E}[f \mid Y] = \mathbb{E}[f]$$

One irrelevant feature isn't a big deal; what we're worried about is when irrelevant features *outnumber* useful ones!

- Decision trees (not too deep)?
Somewhat protected, but beware spurious correlations!
- K -nearest neighbors?

Irrelevant Features

$$\mathbb{E}[f \mid Y] = \mathbb{E}[f]$$

One irrelevant feature isn't a big deal; what we're worried about is when irrelevant features *outnumber* useful ones!

- Decision trees (not too deep)?
Somewhat protected, but beware spurious correlations!
- K -nearest neighbors? 😞

Irrelevant Features

$$\mathbb{E}[f \mid Y] = \mathbb{E}[f]$$

One irrelevant feature isn't a big deal; what we're worried about is when irrelevant features *outnumber* useful ones!

- Decision trees (not too deep)?
Somewhat protected, but beware spurious correlations!
- K -nearest neighbors? 😐
- Perceptron?

Irrelevant Features

$$\mathbb{E}[f \mid Y] = \mathbb{E}[f]$$

One irrelevant feature isn't a big deal; what we're worried about is when irrelevant features *outnumber* useful ones!

- Decision trees (not too deep)?
Somewhat protected, but beware spurious correlations!
- K -nearest neighbors? 😐
- Perceptron? 😊

Irrelevant Features

$$\mathbb{E}[f \mid Y] = \mathbb{E}[f]$$

One irrelevant feature isn't a big deal; what we're worried about is when irrelevant features *outnumber* useful ones!

- Decision trees (not too deep)?
Somewhat protected, but beware spurious correlations!
- K -nearest neighbors? 😞
- Perceptron? 😊

What about *redundant* features ϕ_j and $\phi_{j'}$ such that $\phi_j \approx \phi_{j'}$?

Technique: Feature Pruning

If a binary feature is present in too small or too large a fraction of D , remove it.

Technique: Feature Pruning

If a binary feature is present in too small or too large a fraction of D , remove it.

Example: $\phi(x) = \llbracket \text{the word } \textit{the} \text{ occurs in document } x \rrbracket$

Technique: Feature Pruning

If a binary feature is present in too small or too large a fraction of D , remove it.

Example: $\phi(x) = \llbracket \text{the word } \textit{the} \text{ occurs in document } x \rrbracket$

Generalization: if a feature has variance (in D) **lower** than some threshold value, remove it.

$$\text{sample_mean}(\phi; D) = \frac{1}{N} \sum_{n=1}^N \phi(x_n) \quad (\text{call it "}\bar{\phi}\text{"})$$

$$\text{sample_variance}(\phi; D) = \frac{1}{N-1} \sum_{n=1}^N (\phi(x_n) - \bar{\phi})^2 \quad (\text{call it "Var}(\phi)\text{"})$$

Technique: Feature Normalization

Center a feature:

$$\phi(x) \rightarrow \phi(x) - \bar{\phi}$$

(This was a required step for principal components analysis!)

Scale a feature. Two choices:

$$\phi(x) \rightarrow \frac{\phi(x)}{\sqrt{\text{Var}(\phi)}} \quad \text{“variance scaling”}$$

$$\phi(x) \rightarrow \frac{\phi(x)}{\max_n |\phi(x_n)|} \quad \text{“absolute scaling”}$$

Technique: Feature Normalization

Center a feature:

$$\phi(x) \rightarrow \phi(x) - \bar{\phi}$$

(This was a required step for principal components analysis!)

Scale a feature. Two choices:

$$\phi(x) \rightarrow \frac{\phi(x)}{\sqrt{\text{Var}(\phi)}} \quad \text{“variance scaling”}$$

$$\phi(x) \rightarrow \frac{\phi(x)}{\max_n |\phi(x_n)|} \quad \text{“absolute scaling”}$$

Remember that you'll need to normalize test data before you test!

Technique: Example Normalization

We have been talking about normalizing columns.

We can also normalize rows. l_2 normalization is commonly used for bag of words.

$$\mathbf{x} = \frac{\mathbf{x}}{||\mathbf{x}||_2} = \frac{\mathbf{x}}{\sqrt{\sum_j \mathbf{x}[j]^2}}$$

Techniques: Creating New Features

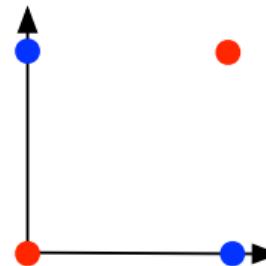
1. Consider two binary features, ϕ_j and $\phi_{j'}$. A new *conjunction* feature can be defined by:

$$\phi_{j \wedge j'}(x) = \phi_j(x) \wedge \phi_{j'}(x)$$

Techniques: Creating New Features

1. Consider two binary features, ϕ_j and $\phi_{j'}$. A new *conjunction* feature can be defined by:

$$\phi_{j \wedge j'}(x) = \phi_j(x) \wedge \phi_{j'}(x)$$

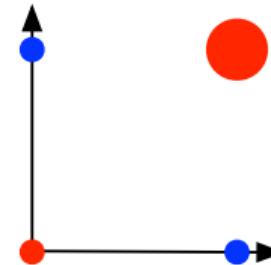


The classic “xor” problem: these points are *not* linearly separable.

Techniques: Creating New Features

1. Consider two binary features, ϕ_j and $\phi_{j'}$. A new *conjunction* feature can be defined by:

$$\phi_{j \wedge j'}(x) = \phi_j(x) \wedge \phi_{j'}(x)$$

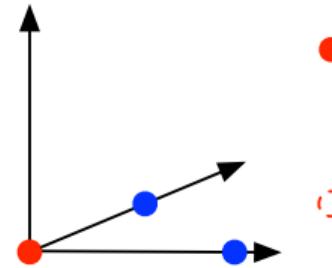


Define $x[3] = x[1] \wedge x[2]$.

Techniques: Creating New Features

1. Consider two binary features, ϕ_j and $\phi_{j'}$. A new *conjunction* feature can be defined by:

$$\phi_{j \wedge j'}(x) = \phi_j(x) \wedge \phi_{j'}(x)$$

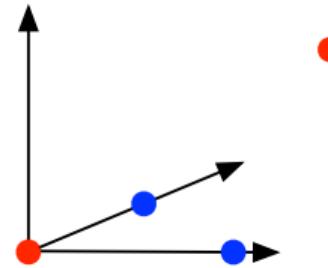


Rotating the view.

Techniques: Creating New Features

1. Consider two binary features, ϕ_j and $\phi_{j'}$. A new *conjunction* feature can be defined by:

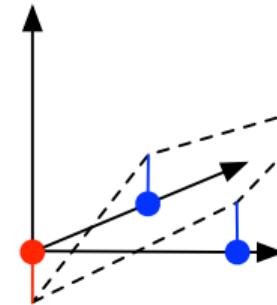
$$\phi_{j \wedge j'}(x) = \phi_j(x) \wedge \phi_{j'}(x)$$



Techniques: Creating New Features

1. Consider two binary features, ϕ_j and $\phi_{j'}$. A new *conjunction* feature can be defined by:

$$\phi_{j \wedge j'}(x) = \phi_j(x) \wedge \phi_{j'}(x)$$



$$2 \cdot \mathbf{x}[1] + 2 \cdot \mathbf{x}[2] - 4 \cdot \mathbf{x}[3] - 1 = 0$$

Techniques: Creating New Features

1. Consider two binary features, ϕ_j and $\phi_{j'}$. A new *conjunction* feature can be defined by:

$$\phi_{j \wedge j'}(x) = \phi_j(x) \wedge \phi_{j'}(x)$$

Generalization: take the *product* of two features.

Techniques: Creating New Features

1. Consider two binary features, ϕ_j and $\phi_{j'}$. A new *conjunction* feature can be defined by:

$$\phi_{j \wedge j'}(x) = \phi_j(x) \wedge \phi_{j'}(x)$$

Generalization: take the *product* of two features.

2. Even more generally, we can create conjunctions (or products) using as many features as we'd like.

Techniques: Creating New Features

1. Consider two binary features, ϕ_j and $\phi_{j'}$. A new *conjunction* feature can be defined by:

$$\phi_{j \wedge j'}(x) = \phi_j(x) \wedge \phi_{j'}(x)$$

- Generalization: take the *product* of two features.
2. Even more generally, we can create conjunctions (or products) using as many features as we'd like.

This is one view of what decision trees are doing!

- Every leaf's path (from root) is a conjunction feature.
- Why not build decision trees, extract the features and toss them into the perceptron?

Techniques: Creating New Features

1. Consider two binary features, ϕ_j and $\phi_{j'}$. A new *conjunction* feature can be defined by:

$$\phi_{j \wedge j'}(x) = \phi_j(x) \wedge \phi_{j'}(x)$$

Generalization: take the *product* of two features.

2. Even more generally, we can create conjunctions (or products) using as many features as we'd like.

This is one view of what decision trees are doing!

- Every leaf's path (from root) is a conjunction feature.
- Why not build decision trees, extract the features and toss them into the perceptron?

3. Transformations on features can be useful. For example,

$$\phi(x) \rightarrow \text{sign}(\phi(x)) \cdot \log(1 + |\phi(x)|)$$

Techniques: Creating New Features

Remember that adding features does not always bring benefits.

You could be just bring irrelevant, redundant, or features that make linearly separable datasets not linearly separable.

A more realistic but easy example

Given the following data about the locations of two cities, predict whether it is possible to drive between these two cities.

City 1 lat.	City 1 long.	City 2 lat.	City 2 long.	drivable
123.24	46.71	121.33	47.34	Yes
123.24	56.91	121.33	55.23	Yes
123.24	46.71	121.33	55.34	No
123.24	46.71	130.99	47.34	No

Feature engineering

Features represent the food of machine learning.

Feature engineering

Features represent the food of machine learning.

Garbage in, garbage out.

Feature engineering

Features represent the food of machine learning.

Garbage in, garbage out.

- Pruning
- Normalization
- Creating new features

Feature engineering

Features represent the food of machine learning.

Garbage in, garbage out.

- Pruning
- Normalization
- Creating new features

In practice, feature engineering requires a deep understanding of the problem.