

CPSC 340: Machine Learning and Data Mining

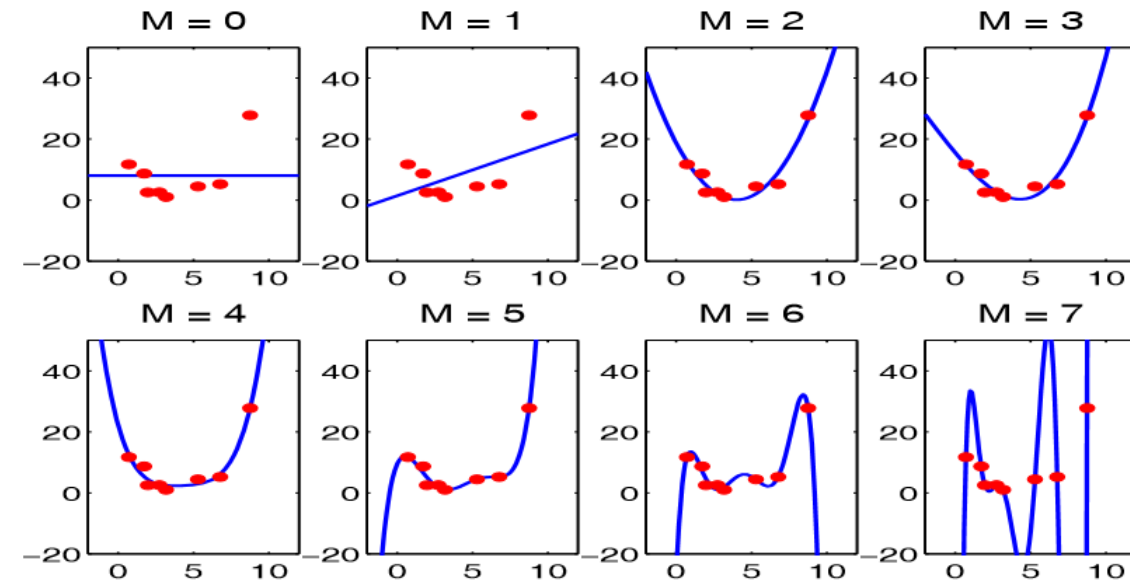
Feature Selection

Fall 2019

Last Time: Finding the “True” Model

- What if y_i really is a polynomial function of x_i ?
 - How can we find the “true” degree ‘p’ of the polynomial?

- Training error does not work:
 - It goes down as ‘p’ goes up.
- Cross-validation may also not work:
 - Tends to overestimate ‘p’.
 - Due to optimization bias.



For example, imagine that the true model is $y_i = 2x_i^2 - 5 + (\text{noise})$

We might choose $d=3$ and a model like $\hat{y}_i = 0.001x_i^3 + 2x_i^2 - 5$,

since it might get a slightly smaller validation error.

Last Time: Complexity Penalties

- We discussed putting a **penalty on the model complexity**.
 - Want to **fit the data and have a simple model**.

Find 'v' and 'p' minimizing:

$$\text{score}(p) = \underbrace{\frac{1}{2} \|Z_p v - y\|^2}_{\text{usual error}} + \underbrace{\lambda K}_{\substack{\text{"strength"} \\ \text{of penalty}}} \quad \rightarrow \text{number of "degrees of freedom"} \\ \text{(K=p+1 for degree-P polynomial)}$$

- “To increase the degrees of freedom by one, need to decrease error by λ ”.
- Prefers smaller degrees of freedom, if errors are similar.
 - **Can't optimize this using gradient descent**, since it's discontinuous in 'p'.
 - Need to search over values of 'p'.

Bayesian Information Criterion (BIC)

- A disadvantage of these methods:
 - Still prefers a larger 'p' as 'n' grows.
- Solution: make λ depend on 'n'.
- For example, the Bayesian information criterion (BIC) uses:
$$\lambda = \frac{1}{2} \log(n)$$
- BIC penalizes a bit more than AIC for large 'n'.
 - As 'n' goes to ∞ , recovers “true” model (“consistent” for model selection).
- In practice, we usually just try a bunch of different λ values.
 - Picking λ is like picking 'k' in k-means.

Discussion of other Scores for Model Selection

- There are many **other scores**:
 - Elbow method (corresponds to specific choice of λ).
 - You could also use BIC for choosing 'k' in k-means.
 - Methods based on validation error.
 - “Take smallest ‘p’ within one standard error of minimum cross-validation error”.
 - Minimum description length.
 - Risk inflation criterion.
 - False discovery rate.
 - **Marginal likelihood** (CPSC 540).
- These can adapted to use the L1-norm and other errors.

Motivation: Discovering Food Allergies

- Recall the food allergy example:

Egg	Milk	Fish	Wheat	Shellfish	Peanuts	...		Sick?
0	0.7	0	0.3	0	0		→	1
0.3	0.7	0	0.6	0	0.01		→	1
0	0	0	0.8	0	0		→	0
0.3	0.7	1.2	0	0.10	0.01		→	1

- Instead of predicting “sick”, we want to do **feature selection**:
 - Which foods are “relevant” for predicting “sick”.

Feature Selection

- General **feature selection** problem:

A hand-drawn diagram illustrating the feature selection problem. It shows a matrix $X = \begin{bmatrix} & & & \end{bmatrix}$ where the second column is circled in blue. An arrow points from the text 'feature j' to this circled column. To the right of X is a vector $y = \begin{bmatrix} \\ \\ \end{bmatrix}$.

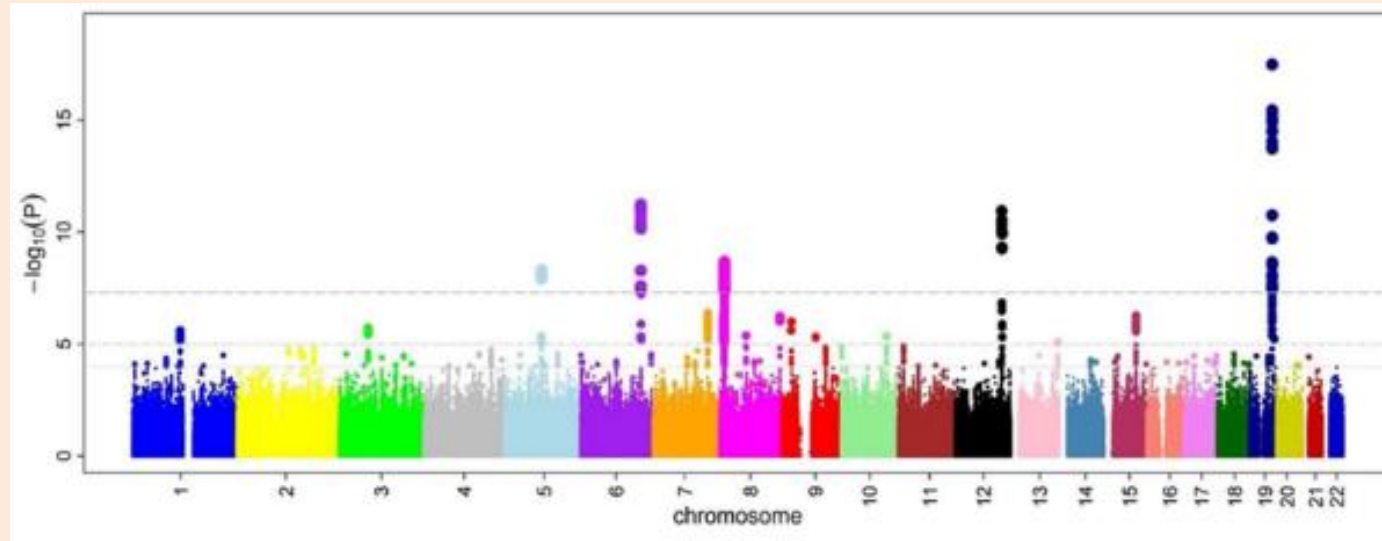
- Find the features (columns) of 'X' that are important for predicting 'y'.
 - “What are the relevant factors?”
 - “Which basis functions should I use among these choices?”
 - “What types of new data should I collect?”
 - “How can I speed up computation?”
- One of most important problems in ML/statistics, **but very messy**.
 - For now, we'll say a feature is “relevant” if it **helps predict y_i from x_i** .

“Association” Approach

- A simple/common way to do feature selection:
 - For each feature ‘j’, compute correlation between feature values x^j and ‘y’.
 - Say that ‘j’ is relevant if correlation is above 0.9 or below -0.9.
- Turns feature selection into hypothesis testing for each feature.
 - There are many other measures of “dependence” ([Wikipedia](#)).
- Usually gives unsatisfactory results as it ignores variable interactions:
 - Includes irrelevant variables: “Taco Tuesdays”.
 - If tacos make you sick, and you often eat tacos on Tuesdays, it will say “Tuesday” is relevant.
 - Excludes relevant variables: “Diet Coke + Mentos Eruption”.
 - Diet coke and Mentos don’t make you sick on their own, but *together* they make you sick.

Genome-Wide Association Studies

- Genome-wide association studies:
 - Measure if there exists a dependency between each individual “single-nucleotide polymorphism” in the genome and a particular disease.



- Has identified thousands of genes “associated” with diseases.
 - But *by design* this has a **huge numbers of false positives** (and many false negatives).

“Regression Weight” Approach

- A simple/common approach to feature selection:
 - Fit regression weights ‘w’ based on **all** features (maybe with least squares).
 - Take all features ‘j’ where **weight $|w_j|$ is greater than a threshold**.
- For example: you fit a least squares model with 5 features and get:

$$w = \begin{bmatrix} 0.01 \\ -0.1 \\ 10 \\ -3 \\ 0.001 \end{bmatrix}$$

- Feature 3 looks the most relevant.
- Feature 4 also looks relevant.
- Feature 5 seems irrelevant.

“Regression Weight” Approach

- A simple/common approach to feature selection:
 - Fit regression weights ‘ w ’ based on **all** features (maybe with least squares).
 - Take all features ‘ j ’ where **weight $|w_j|$ is greater than a threshold.**
- This could recognize that “Tuesday” is irrelevant.
 - If you get enough data, and you sometimes eat tacos on other days.
(And the relationship is actually linear.)
- This could recognize that “Diet Coke” and “Mentos” are relevant.
 - Assuming this combination occurs enough times in the data.

“Regression Weight” Approach

- A simple/common approach to feature selection:
 - Fit regression weights ‘w’ based on **all** features (maybe with least squares).
 - Take all features ‘j’ where **weight $|w_j|$ is greater than a threshold**.

- Has **major problems with collinearity**:

- If the “Tuesday” variable always equals the “taco” variable, it **could say that Tuesdays are relevant but tacos are not**.

$$\hat{y}_i = w_1 * \text{taco} + w_2 * \text{Tuesday} = 0 * \text{taco} + (w_1 + w_2) * \text{Tuesday}$$

- If you have two copies of an irrelevant feature, it **could take both irrelevant copies**.

$$\hat{y}_i = 0 * \text{irrelevant} + 0 * \text{irrelevant} = 10000 * \text{irrelevant} + (-10000) * \text{irrelevant}$$

(pause)

Search and Score Methods

- Most common feature selection framework is **search and score**:
 1. Define **score function $f(S)$** that measures quality of a set of features 'S'.
 2. Now **search** for the variables 'S' with the best score.
- Example with 3 features:
 - Compute “score” of using feature 1.
 - Compute “score” of using feature 2.
 - Compute “score” of using feature 3.
 - Compute “score” of using features {1,2}.
 - Compute “score” of using features {1,3}.
 - Compute “score” of using features {2,3}.
 - Compute “score” of using features {1,2,3}.
 - Compute “score” of using features {}.
 - Return the set of features 'S' with the best “score”.

Which Score Function?

- The **score can't be the training error**.
 - Training error goes down as you add features, so will **select all features**.
- A more logical score is the **validation error**.
 - “**Find the set of features that gives the lowest validation error.**”
 - To minimize test error, this is what we want.
- But there are problems due to the **large number of sets** of variables:
 - If we have 'd' variables, there are **2^d sets** of variables.
 - **Optimization bias** is high: we're optimizing over 2^d models (not 10).
 - Prone to **false positives**: irrelevant variables will sometimes help by chance.

“Number of Features” Penalties

- To reduce false positives, we can again use complexity penalties:

$$\text{score}(S) = \frac{1}{2} \sum_{i=1}^n (w_S^T x_{iS} - y_i)^2 + \text{size}(S)$$

- E.g., we could use **squared error** and **number of non-zeroes**.
- We’re using ‘ x_{iS} ’ as the features ‘ S ’ of example x_i .
- If two ‘ S ’ have similar error, this **prefers the smaller set**.
 - It **prefers removing feature 3 instead of having $w_3 = 0.00001$** .
- Instead of “ $\text{size}(S)$ ”, we usually write this using the “L0-norm”...

L0-Norm and “Number of Features We Use”

- In linear models, **setting $w_j = 0$ is the same as removing feature ‘j’**:

$$\begin{aligned}\hat{y}_i &= w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i3} + \dots + w_d x_{id} \\ &\quad \quad \quad \underbrace{\hspace{1.5cm}}_{\text{set } w_2 = 0} \\ \hat{y}_i &= w_1 x_{i1} + \underbrace{0}_{\text{ignore } x_{i2}} + w_3 x_{i3} + \dots + w_d x_{id}\end{aligned}$$

- The L0 “norm” is the number of non-zero values ($\|w\|_0 = \text{size}(S)$).

$$\text{If } w = \begin{bmatrix} 1 \\ 0 \\ 2 \\ 0 \\ 3 \end{bmatrix} \text{ then } \|w\|_0 = 3 \qquad \text{If } w = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \text{ then } \|w\|_0 = 0.$$

- Not actually a true norm.
- If ‘w’ has a small L0-norm, then it doesn’t use many features.

L0-penalty: optimization

- L0-norm penalty for feature selection:

$$f(w) = \underbrace{\frac{1}{2} \|Xw - y\|^2}_{\text{training error}} + \underbrace{\lambda \|w\|_0}_{\text{degrees of freedom 'k'}}$$

- Suppose we want to use this to evaluate the features $S = \{1, 2\}$:
 - First fit the 'w' just using features 1 and 2.
 - Now compute the training error with this 'w' and features 1 and 2.
 - Add $\lambda * 2$ to the training error to get the score.
- We repeat this with other choices of 'S' to find the “best” features.

L0-penalty: interpretation

- L0-norm penalty for feature selection:

$$f(w) = \frac{1}{2} \|Xw - y\|^2 + \lambda \|w\|_0$$

- Balances between training error and number of features we use.
 - With $\lambda=0$, we get least squares with all features.
 - With $\lambda=\infty$, we must set $w=0$ and not use any features.
 - With other λ , balances between training error and number of non-zeroes.
 - Larger λ puts more emphasis on having zeroes in 'w' (more feature selection).
 - Different values give AIC, BIC, and so on.

Forward Selection (Greedy Search Heuristic)

- In **search and score**, it's also just **hard to search for the best 'S'**.
 - There are **2^d possible sets**.
- A common greedy search procedure is **forward selection**:
 1. Compute score if we use no features.
 2. Try adding "taco", "milk", "egg", and so on (computing score of each)
 3. Add "milk" because it got the best score.
 4. Try {milk, taco}, {milk, egg}, and so on (computing score of each variable with milk)
 5. Add "egg" because it got the best score combined with "milk"
 6. Try {milk, egg, taco}, {milk, egg, pizza}, and so on...

Forward Selection (Greedy Search Heuristic)

- **Forward selection** algorithm for variable selection:
 1. Start with an **empty set** of features, $S = []$.
 2. For each possible feature 'j':
 - **Compute scores of features in 'S' combined with feature 'j'.**
 3. Find the 'j' that has the highest score when added to 'S'.
 4. Check if $\{S \cup j\}$ improves on the best score found so far.
 5. Add 'j' to 'S' and go back to Step 2.
 - A variation is to **stop if no 'j' improves the score** over just using 'S'.
- **Not guaranteed to find the best** set, but **reduces many problems**:
 - Considers $O(d^2)$ models: cheaper, overfits less, has fewer false positives.

Backward Selection and RFE

- **Forward selection** often works better than naïve methods.
- A related method is **backward selection**:
 - Start with all features, compute score after removing each feature, remove the one that improves the score the most.
- If you consider adding or removing features, it's called **stagewise**.
- **Stochastic local search** is a class of fancier methods.
 - Simulated annealing, genetic algorithms, ant colony optimization, etc.
- **Recursive feature elimination** is another related method:
 - Fit parameters of a regression model.
 - Prune features with small regression weights.
 - Repeat.

(pause)

Is “Relevance” Clearly Defined?

- Consider a supervised classification task:

gender	mom	dad
F	1	0
M	0	1
F	0	0
F	1	1

SNP
1
0
0
1

- Predict whether someone has particular genetic variation (SNP).
 - Location of mutation is in “mitochondrial” DNA.
 - “You almost always have the same value as your mom”.
 - For simplicity we’ll assume 1950s-style gender and parentage.

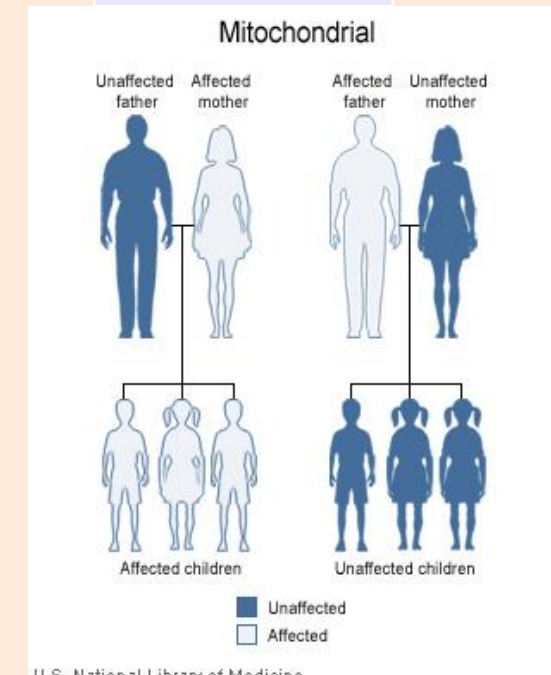
Is “Relevance” Clearly Defined?

- Consider a supervised classification task:

gender	mom	dad
F	1	0
M	0	1
F	0	0
F	1	1

SNP
1
0
0
1

- True model:
 - (SNP = mom) with very high probability.
 - (SNP != mom) with some very low probability.
- What are the “relevant” features for this problem?
 - Mom is relevant and {gender, dad} are not relevant.



Is “Relevance” Clearly Defined?

- What if “mom” feature is repeated?

gender	mom	dad	mom2
F	1	0	1
M	0	1	0
F	0	0	0
F	1	1	1

SNP
1
0
0
1

- Are “mom” and “mom2” relevant?
 - Should we pick them both?
 - Should we pick one because it predicts the other?
- General problem (“dependence”, “collinearity” for linear models):
 - If **features can be predicted from features**, don’t know one(s) to pick.

*Neither of these
is “correct”, but
not picking either
is incorrect.*

Is “Relevance” Clearly Defined?

- What if we add (maternal) “grandma”?

gender	mom	dad	grandma
F	1	0	1
M	0	1	0
F	0	0	0
F	1	1	1

SNP
1
0
0
1

- Is “grandma” relevant?
 - You can predict SNP very accurately from “grandma” alone.
 - But “grandma” is irrelevant if I know “mom”.
- General problem (**conditional independence**):
 - “Relevant” features may be **irrelevant given other features**.

Is “Relevance” Clearly Defined?

- What if we don’t know “mom”?

gender	grandma	dad
F	1	0
M	0	1
F	0	0
F	1	1

SNP
1
0
0
1

- Now is “grandma” is relevant?
 - Without “mom” variable, using “grandma” is the best you can do.
- General problem (“taco Tuesday”):
 - Features can be **relevant due to missing information**.

Is “Relevance” Clearly Defined?

- What if we don’t know “mom” or “grandma”?

gender	dad
F	0
M	1
F	0
F	1

SNP
1
0
0
1

- Now there are no relevant variables, right?
 - But “dad” and “mom” must have some common maternal ancestor.
 - “Mitochondrial Eve” estimated to be ~200,000 years ago.
- General problem (effect size):
 - “Relevant” features may have small effects.

Is “Relevance” Clearly Defined?

- What if we don’t know “mom” or “grandma”?

gender	dad
F	0
M	1
F	0
F	1

SNP
1
0
0
1

- Now there are no relevant variables, right?
 - What if “mom” likes “dad” because he has the same SNP as her?
- General problem (**confounding**):
 - Hidden effects can **make “irrelevant” variables “relevant”**.

Is “Relevance” Clearly Defined?

- What if we add “sibling”?

gender	dad	sibling
F	0	1
M	1	0
F	0	0
F	1	1

SNP
1
0
0
1

- Sibling is “relevant” for predicting SNP, but it’s not the cause.
- General problem (non-causality or reverse **causality**):
 - A “relevant” feature **may not be causal**, or may be an effect of label.

Is “Relevance” Clearly Defined?

- What if don't have “mom” but we have “baby”?

gender	dad	baby
F	0	1
M	1	1
F	0	0
F	1	1

SNP
1
0
0
1

- “Baby” is relevant when (gender == F).
 - “Baby” is relevant (though causality is reversed).
 - Is “gender” relevant?
 - If we want to find relevant causal factors, “gender” is not relevant.
 - If we want to predict SNP, “gender” is relevant.
- General problem (**context-specific relevance**):
 - Adding a feature can **make an “irrelevant” feature “relevant”**.

Is “Relevance” Clearly Defined?

- Warnings about feature selection:
 - A feature is **only “relevant” in the context of available features.**
 - Adding/removing features can make features relevant/irrelevant.
 - Confounding factors can **make “irrelevant” variables the most “relevant”.**
 - If features can be predicted from features, **you can’t know which to pick.**
 - Collinearity is a special case of “dependence” (which may be non-linear).
 - A “relevant” feature may have a **tiny effect.**
 - “Relevance” for prediction does **not imply a causal relationship.**

Is this hopeless?

- We often want to do feature selection we so have to try!
- Different methods are affected by problems in different ways.
 - We'll ignore causality and confounding issues (bonus slides).
- These “problems” don't have right answers but have **wrong answers**:
 - **Variable dependence** (“mom” and “mom2” have same information).
 - But should take at least one.
 - **Conditional independence** (“grandma” is irrelevant given “mom”).
- These “problems” have **application-specific answers**:
 - **Tiny effects**.
 - **Context-specific relevance** (is “gender” relevant if given “baby”?).

Rough Guide to Feature Selection

Method\Issue	Dependence	Conditional Independence	Tiny effects	Context-Specific Relevance
Association (e.g., measure correlation between features 'j' and 'y')	Ok (takes "mom" and "mom2")	Bad (takes "grandma", "great-grandma", etc.)	Ignores	Bad (misses features that must interact, "gender" irrelevant given "baby")

Rough Guide to Feature Selection

Method\Issue	Dependence	Conditional Independence	Tiny effects	Context-Specific Relevance
Association (e.g., measure correlation between features 'j' and 'y')	Ok (takes "mom" and "mom2")	Bad (takes "grandma", "great-grandma", etc.)	Ignores	Bad (misses features that must interact, "gender" irrelevant given "baby")
Regression Weight (fit least squares, take biggest $ w_j $)	Bad (can take irrelevant but collinear, can take none of "mom1-3")	Ok (takes "mom" not "grandma", if linear and 'n' large.	Ignores (unless collinear)	Ok (if linear, "gender" relevant give "baby")

Rough Guide to Feature Selection

Method\Issue	Dependence	Conditional Independence	Tiny effects	Context-Specific Relevance
Association (e.g., measure correlation between features 'j' and 'y')	Ok (takes "mom" and "mom2")	Bad (takes "grandma", "great-grandma", etc.)	Ignores	Bad (misses features that must interact, "gender" irrelevant given "baby")
Regression Weight (fit least squares, take biggest $ w_j $)	Bad (can take irrelevant but collinear, can take none of "mom1-3")	Ok (takes "mom" not "grandma", if linear and 'n' large.	Ignores (unless collinear)	Ok (if linear, "gender" relevant give "baby")
Search and Score w/ Validation Error	Ok (takes at least one of "mom" and "mom2")	Bad (takes "grandma", "great-grandma", etc.)	Allows	Ok (“gender” relevant given “baby”)

Rough Guide to Feature Selection

Method\Issue	Dependence	Conditional Independence	Tiny effects	Context-Specific Relevance
Association (e.g., measure correlation between features 'j' and 'y')	Ok (takes "mom" and "mom2")	Bad (takes "grandma", "great-grandma", etc.)	Ignores	Bad (misses features that must interact, "gender" irrelevant given "baby")
Regression Weight (fit least squares, take biggest $ w_j $)	Bad (can take irrelevant but collinear, can take none of "mom1-3")	Ok (takes "mom" not "grandma", if linear and 'n' large.	Ignores (unless collinear)	Ok (if linear, "gender" relevant give "baby")
Search and Score w/ Validation Error	Ok (takes at least one of "mom" and "mom2")	Bad (takes "grandma", "great-grandma", etc.)	Allows (many false positives)	Ok (“gender” relevant given “baby”)
Search and Score w/ L0-norm	Ok (takes exactly one of "mom" and "mom2")	Ok (takes "mom" not grandma if linear-ish).	Ignores (even if collinear)	Ok (“gender” relevant given “baby”)

Summary

- **Feature selection** is task of choosing the “relevant” features.
 - Obvious simple approaches have obvious simple problems.
- **Search and score**: find features that optimize some score.
 - **L0-norm penalties** are the most common scores.
 - **Forward selection** is a heuristic to search over a smaller set of features.
- “**Relevance**” is really hard to define.
- Next time: getting a good test error even with irrelevant features.

Mallow's Cp

- Older than AIC and BIC is **Mallow's Cp**:

$$f(w) = \frac{\|Xw - y\|^2}{\frac{1}{n} \|X\hat{w} - y\|^2} - n + 2\|w\|_0$$

least squares weights if we used all features.

- Minimizing this score is equivalent to L0-regularization:

$$f(w) = \frac{1}{2} \|Xw - y\|^2 + \lambda \|w\|_0$$

$$\text{with } \lambda = \frac{\|X\hat{w} - y\|^2}{n}$$

- So again, viewing λ as hyper-parameter, this score is special case.

Adjusted R^2

- Older than AIC and BIC and Mallow's C_p is **adjusted R^2** :

$$f(w) = 1 - (1 - R^2) \frac{n-1}{n - \|w\|_0 - 1} \quad \text{where} \quad R^2 = 1 - \frac{\|Xw - y\|^2}{\|X\hat{w} - y\|^2}$$

- Maximizing this score is equivalent to L0-regularization:

$$= \frac{1}{2} \|Xw - y\|^2 + \lambda \|w\|_0$$

$$\text{with } \lambda = \frac{\|X\hat{w} - y\|^2}{2(n-1)}$$

- So again, viewing λ as hyper-parameter, this score is special case.

ANOVA

- Some people also like to compute this “ANOVA” quantity:

$$f(w) = \frac{\|Xw - \bar{y}\|^2}{\|y - \bar{y}\|^2}$$

↑ mean of y_i values repeated n times

- This is based on the decomposition of “total squared error” as:

$$\underbrace{\|y - \bar{y}\|^2}_{\text{“total” error}} = \underbrace{\|Xw - \bar{y}\|^2}_{\text{“explained” error}} + \underbrace{\|Xw - y\|^2}_{\text{“residual” (usual) error}}$$

- Notice that “explained error” goes up as our usual (“residual”) error goes down.
- Trying to find the ‘k’ features that maximize ‘f’ (“explain the most variance”) is equivalent to L0-regularization with a particular λ (so another special case).

Information Criteria with Noise Variance

- We defined AIC/BIC for feature selection in least squares as:

$$f(w) = \frac{1}{2} \|Xw - y\|^2 + \lambda \|w\|_0$$

- The first term comes from assuming $y_i = w^T x_i + \varepsilon$, where ε comes from a normal distribution with a variance of 1.
 - We'll discuss why when we discuss MLE and MAP estimation.
 - If you aren't doing least squares, replace first term by "log-likelihood".
- If you treat variance as a parameter, then after some manipulation:

$$f(w) = \frac{n}{2} \log(\|Xw - y\|^2) + \lambda \|w\|_0$$

- However, this is again equivalent to just changing λ .

Complexity Penalties for Other Models

- Scores like AIC and BIC can also be used in other contexts:
 - When fitting a decision tree, only split a node if it improves BIC.
 - This makes sense if we're looking for the “true tree”, or maybe just a simple/interpretable tree that performs well.
- In these cases we replace “L0-norm” with “degrees of freedom”.
 - In linear models fit with least squares, degrees of freedom is number of non-zeroes.
 - Unfortunately, it is not always easy to measure “degrees of freedom”.

Multi-Dimensional Polynomial Basis

- Recall fitting **polynomials** when we only have 1 feature:

$$\hat{y}_i = w_0 + w_1 x_i + w_2 x_i^2$$

- We can fit these models using a **change of basis**:

$$X = \begin{bmatrix} 0.2 \\ -0.5 \\ 1 \\ 4 \end{bmatrix} \quad Z = \begin{bmatrix} 1 & 0.2 & (0.2)^2 \\ 1 & -0.5 & (-0.5)^2 \\ 1 & 1 & (1)^2 \\ 1 & 4 & (4)^2 \end{bmatrix}$$

- How can we do this when we have a lot of features?

Multi-Dimensional Polynomial Basis

- Polynomial basis for $d=2$ and $p=2$:

$$X = \begin{bmatrix} 0.2 & 0.3 \\ 1 & 0.5 \\ -0.5 & -0.1 \end{bmatrix} \longrightarrow Z = \begin{bmatrix} 1 & 0.2 & 0.3 & (0.2)^2 & (0.3)^2 & (0.2)(0.3) \\ 1 & 1 & 0.5 & (1)^2 & (0.5)^2 & (1)(0.5) \\ 1 & 0.5 & -0.1 & (0.5)^2 & (-0.1)^2 & (-0.5)(-0.1) \end{bmatrix}$$

$\underbrace{\hspace{1.5cm}}_{\text{bias}}$
 $\underbrace{\hspace{1.5cm}}_{x_{i1}}$
 $\underbrace{\hspace{1.5cm}}_{x_{i2}}$
 $\underbrace{\hspace{1.5cm}}_{(x_{i1})^2}$
 $\underbrace{\hspace{1.5cm}}_{(x_{i2})^2}$
 $\underbrace{\hspace{1.5cm}}_{(x_{i1})(x_{i2})}$

- With $d=4$ and $p=3$, the polynomial basis would include:
 - Bias variable and the x_{ij} : $1, x_{i1}, x_{i2}, x_{i3}, x_{i4}$.
 - The x_{ij} squared and cubed: $(x_{i1})^2, (x_{i2})^2, (x_{i3})^2, (x_{i4})^2, (x_{i1})^3, (x_{i2})^3, (x_{i3})^3, (x_{i4})^3$.
 - Two-term interactions: $x_{i1}x_{i2}, x_{i1}x_{i3}, x_{i1}x_{i4}, x_{i2}x_{i3}, x_{i2}x_{i4}, x_{i3}x_{i4}$.
 - Cubic interactions: $x_{i1}x_{i2}x_{i3}, x_{i2}x_{i3}x_{i4}, x_{i1}x_{i3}x_{i4}, x_{i1}x_{i2}x_{i4},$
 $x_{i1}^2x_{i2}, x_{i1}^2x_{i3}, x_{i1}^2x_{i4}, x_{i1}x_{i2}^2, x_{i2}^2x_{i3}, x_{i2}^2x_{i4}, x_{i1}x_{i3}^2, x_{i2}x_{i3}^2, x_{i3}^2x_{i4}, x_{i1}x_{i4}^2, x_{i2}x_{i4}^2, x_{i3}x_{i4}^2$.

Kernel Trick

- If we go to degree $p=5$, we'll have $O(d^5)$ quintic terms:

$$x_{i1}^5, x_{i1}^4 x_{i2}, x_{i1}^4 x_{i3}, \dots, x_{i1}^4 x_{id}, x_{i1}^3 x_{i2}^2, x_{i1}^3 x_{i3}^2, \dots, x_{i1}^3 x_{id}^2, \dots, x_{i2}^5, x_{i2}^4 x_{i3}, \dots, x_{id}^5$$

- For large 'd' and 'p', **storing a polynomial basis is intractable!**
 - 'Z' has $k=O(d^p)$ columns, so it does not fit in memory.
- One way to handle this is to use the “**polynomial kernel**” (later).
 - Computational trick that lets you use all these features efficiently.
- Alternately, choose a subset with **hierarchical forward selection**.

Feature Selection Hierarchy

- Consider a linear models with **higher-order terms**,

$$\hat{y}_i = w_0 + w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i3} + w_{12} x_{i1} x_{i2} + w_{13} x_{i1} x_{i3} + w_{23} x_{i2} x_{i3} + w_{123} x_{i1} x_{i2} x_{i3}$$

- The **number of higher-order terms may be too large**.
 - Can't even compute them all.
 - We need to somehow decide which terms we'll even consider.
- Consider the following **hierarchical constraint**:
 - You only **allow** $w_{12} \neq 0$ if $w_1 \neq 0$ and $w_2 \neq 0$.
 - "Only consider feature interaction if you are using both features already."

Hierarchical Forward Selection

- Hierarchical Forward Selection:
 - Usual forward selection, but **consider interaction terms obeying hierarchy**.
 - Only consider $w_{12} \neq 0$ once $w_1 \neq 0$ and $w_2 \neq 0$.
 - Only allow $w_{123} \neq 0$ once $w_{12} \neq 0$ and $w_{13} \neq 0$ and $w_{23} \neq 0$.
 - Only allow $w_{1234} \neq 0$ once all threeway interactions are present.

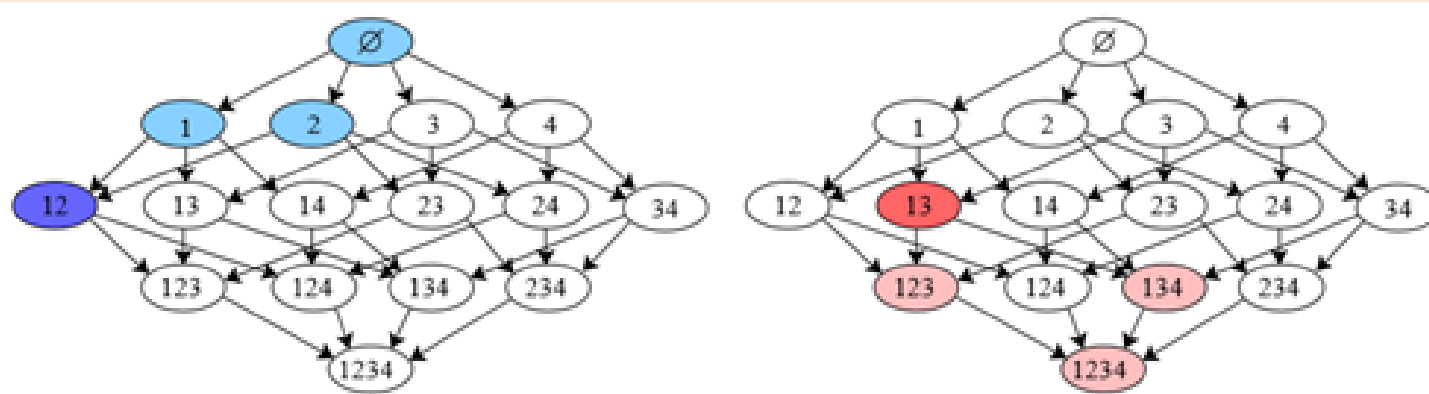
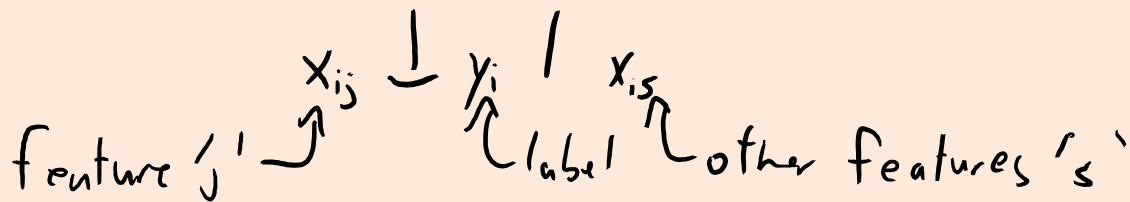


Fig 9: Power set of the set $\{1, \dots, 4\}$: in blue, an authorized set of selected subsets. In red, an example of a group used within the norm (a subset and all of its descendants in the DAG).

Alternative to Search and Score: good old p-values

- **Hypothesis testing** (“constraint-based”) approach:
 - Generalization of the “association” approach to feature selection.
 - Performs a sequence of **conditional independence tests**.



“If I know features in ‘s’
does feature ‘j’ tell me anything
about label?”

- If they are independent (like “ $p < .05$ ”), say that ‘j’ is “irrelevant”.
- Common way to do the tests:
 - “Partial” correlation (numerical data).
 - “Conditional” mutual information (discrete data).

Testing-Based Feature Selection

- Hypothesis testing (“constraint-based”) approach:
- Too many possible tests, “greedy” method is for each ‘j’ do:

First test if $x_{ij} \perp y_i$

If still dependent test $x_{ij} \perp y_i \mid x_{i,s}$ where ‘s’ has one feature

If still dependent test $x_{ij} \perp y_i \mid x_{i,s}$ where ‘s’ now has two features dependence.

Often choose features to minimize dependence.

⋮
If still dependent when ‘s’ includes all other features, declare ‘j’ relevant.

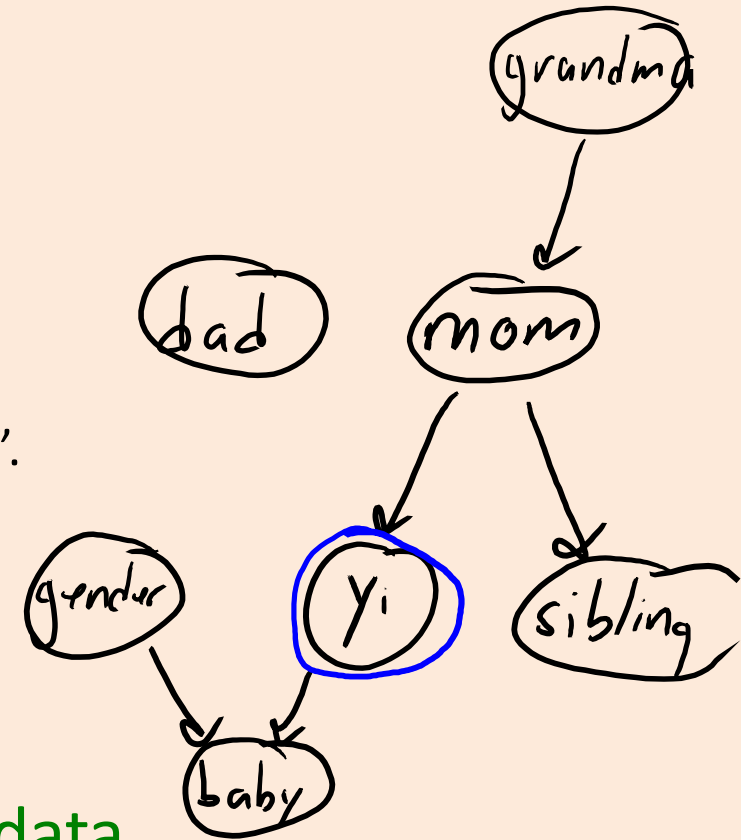
- “Association approach” is the greedy method where you **only do the first test** (subsequent tests remove a lot of false positives).

Hypothesis-Based Feature Selection

- Advantages:
 - Deals with conditional independence.
 - Algorithm can explain why it thinks ‘j’ is irrelevant.
 - Doesn’t necessarily need linearity.
- Disadvantages:
 - Deals badly with exact dependence: doesn’t select “mom” or “mom2” if both present.
 - Usual warning about testing multiple hypotheses:
 - If you test $p < 0.05$ more than 20 times, you’re going to make errors.
 - Greedy approach may be sub-optimal.
- Neither good nor bad:
 - Allows tiny effects.
 - Says “gender” is irrelevant when you know “baby”.
 - This approach is sometimes better for finding relevant factors, not to select features for learning.

Causality

- None of these approaches address **causality or confounding**:
 - “Mom” is the **only relevant direct causal factor**.
 - “Dad” is really irrelevant.
 - “Grandma” is causal but is irrelevant if we know “mom”.
- Other factors can **help prediction but aren’t causal**:
 - “Sibling” is predictive due to **confounding** of effect of same “mom”.
 - “Baby” is predictive due to **reverse causality**.
 - “Gender” is predictive due to **common effect** on “baby”.
- We can sometimes address this using **interventional data...**



Interventional Data

- The difference between **observational** and **interventional** data:
 - If I **see** that my watch says 10:45, class is almost over (**observational**).
 - If I **set** my watch to say 10:45, it doesn't help (**interventional**).
- The **intervention** can help discover causal effects:
 - “Watch” is only predictive of “time” in observational setting (so not causal).
- General idea for **identifying causal effects**:
 - “Force” the variable to take a certain value, then measure the effect.
 - If the dependency remains, there is a causal effect.
 - We “break” connections from reverse causality, common effects, or confounding.

Causality and Dataset Collection

- This has to do with the **way you collect data**:
 - You **can't "look" for variables taking the value** "after the fact".
 - You **need to manipulate the value of the variable**, then watch for changes.
- This is the basis for **randomized control trial** in medicine:
 - Randomly assigning pills "forces" value of "treatment" variable.
 - Include a "control" as a value to prevent placebo effect as confounding.
- See also Simpson's Paradox:
 - <https://www.youtube.com/watch?v=ebEkn-BiW5k>

Structure Learning: Unsupervised Feature Selection

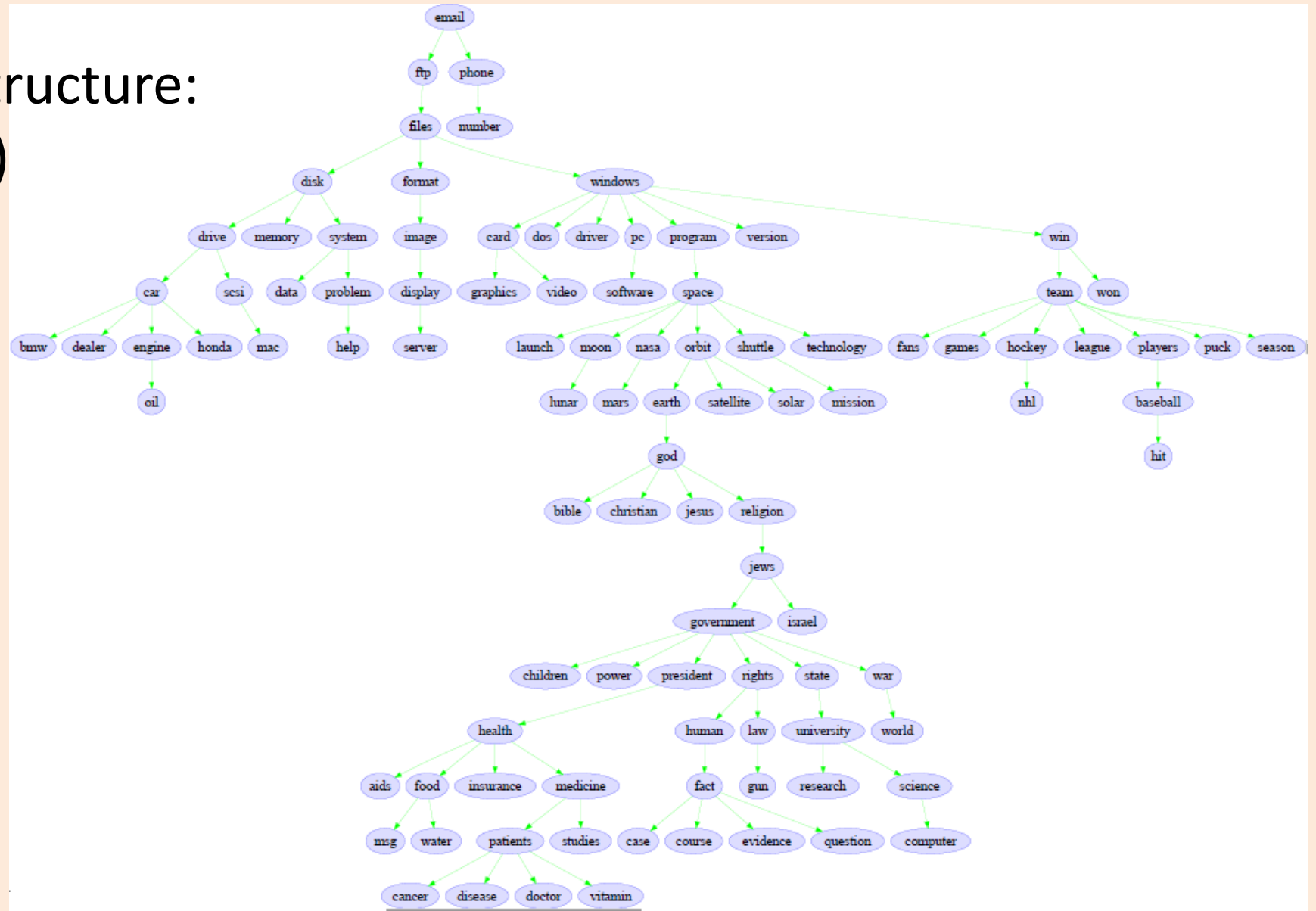
- “News” data: presence of 100 words in 16k newsgroup posts:

car	drive	files	hockey	mac	league	pc	win
0	0	1	0	1	0	1	0
0	0	0	1	0	1	0	1
1	1	0	0	0	0	0	0
0	1	1	0	1	0	0	0
0	0	1	0	0	0	1	1

- Which words are related to each other?
- Problem of **structure learning**: **unsupervised feature selection**.

Structure Learning: Unsupervised Feature Selection

- Optimal tree structure:
(ignore arrows)



Naïve Approach: Association Networks

- A naïve approach to structure learning (“association networks”):
 - For each pair of variables, compute a measure of similarity or dependence.
- Using these n^2 similarity values either:
 - Select all pairs whose similarity is above a threshold.
 - Select the “top k” most similar features to each feature ‘j’.
- Main problems:
 - Usually, most variables are dependent (too many edges).
 - “Sick” is getting connected to “Tuesdays” even if “tacos” are a variable.
 - “True” neighbours may not have the highest dependence.
 - “Sick” might get connected to “Tuesdays” before it gets connected to “milk”.
- (Variation: best tree can be found as minimum spanning tree problem.)

Example: Vancouver Rain Data

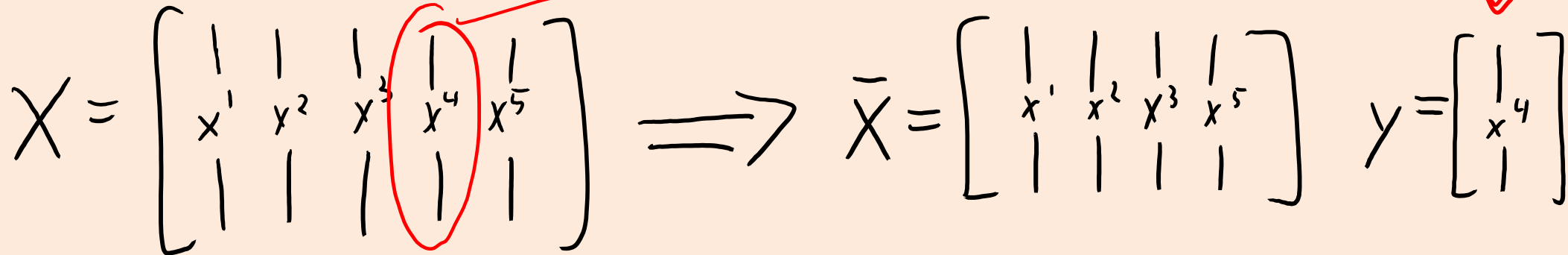
- Consider modeling the “Vancouver rain” dataset.

	Day 1	Day 2	Day 3	Day 4	Day 5	Day 6	Day 7	Day 8	Day 9	...
Month 1	0	0	0	1	1	0	0	1	1	
Month 2	1	0	0	0	0	0	1	0	0	
Month 3	1	1	1	1	1	1	1	1	1	
Month 4	1	1	1	1	0	0	1	1	1	
Month 5	0	0	0	0	1	1	0	0	0	
Month 6	0	1	1	0	0	0	0	1	1	

- The strongest signal in the data is the simple relationship:
 - If it rained yesterday, it's likely to rain today ($> 50\%$ chance that $x^{t-1} = x^t$).
 - But an “association network” might connect all days (all dependent).

Dependency Networks

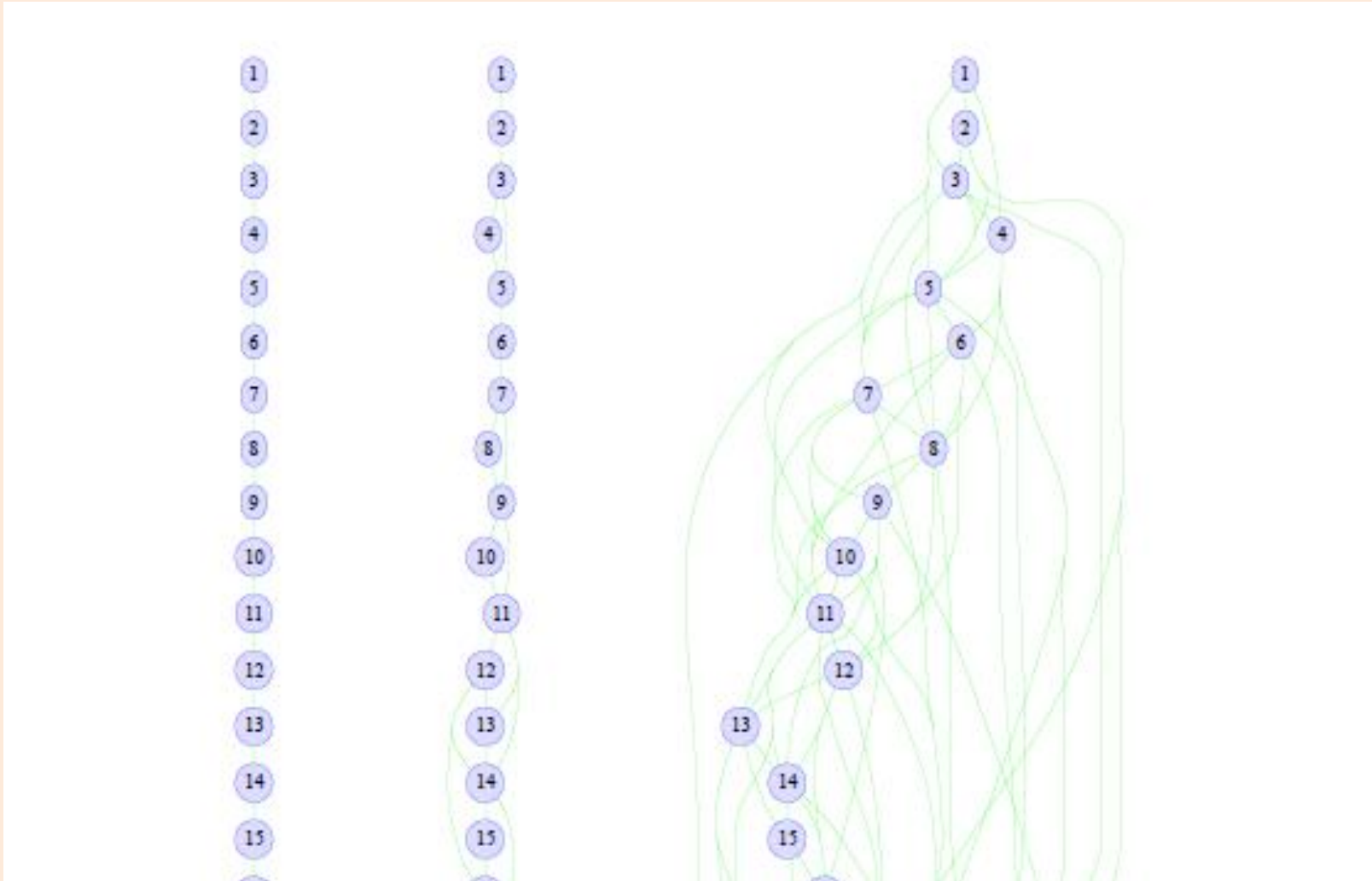
- A better approach is **dependency networks**:
 - For each variable 'j', **make it the target in a supervised learning problem.**


$$X = \begin{bmatrix} | & | & | & | & | \\ x^1 & x^2 & x^3 & x^4 & x^5 \\ | & | & | & | & | \end{bmatrix} \Rightarrow \bar{X} = \begin{bmatrix} | & | & | & | \\ x^1 & x^2 & x^3 & x^5 \\ | & | & | & | \end{bmatrix} \quad y = \begin{bmatrix} | \\ x^4 \\ | \end{bmatrix}$$

- Now we can **use any feature selection method** to choose j's "neighbours".
 - Forward selection, L1-regularization, ensemble methods, etc.
- Can capture **conditional independence**:
 - Might connect "sick" to "tacos", and "tacos" to "Tuesdays" (w/o sick-tacos).

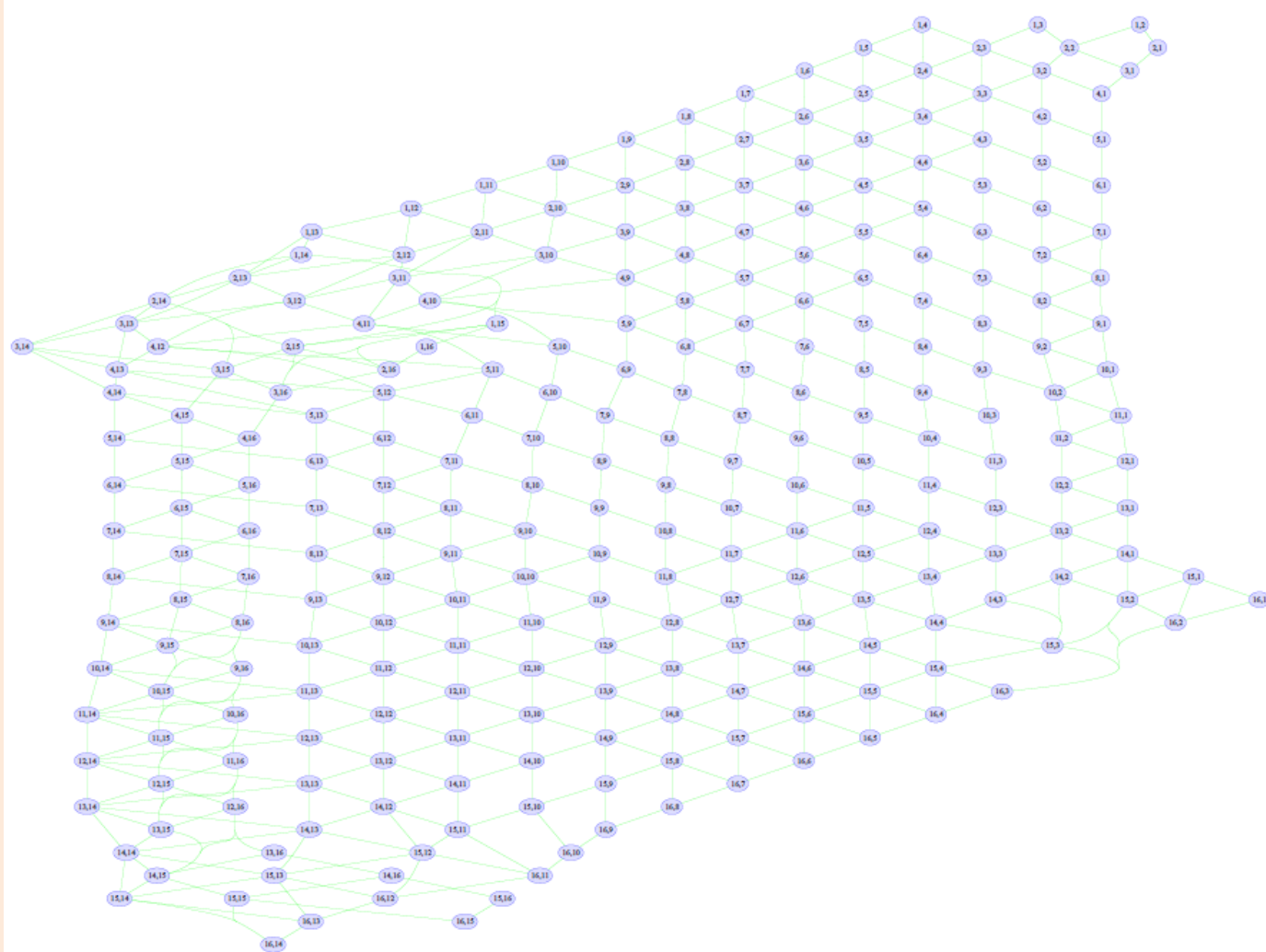
Dependency Networks

- Dependency network fit to Vancouver rain data (different λ values):



Dependency Networks

- Variation on dependency networks on digit image pixels:



Another popular structure learning method is the "PC" algorithm.