# CPSC 340:
# Machine Learning and Data Mining
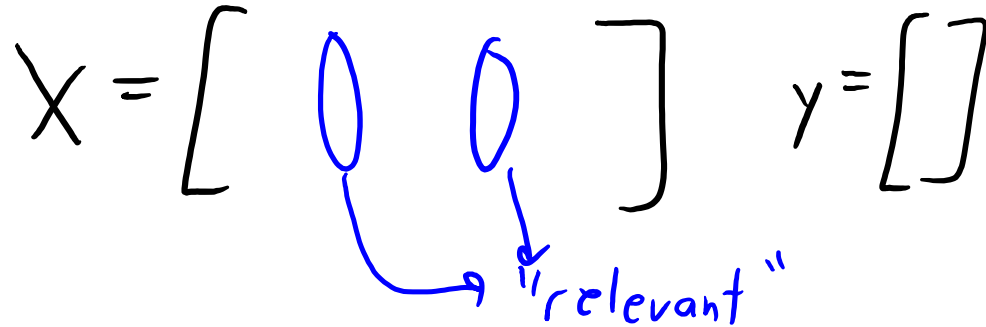
## Regularization

# Admin

- ## Assignment 3
  - Due Friday

- ## Midterm
  - Feb 14 (a week from today) in class
  - If your surname starts with the letters A-G, room DMP 201
  - If your surname starts with the letters H-Z, room DMP 110 (this room)

# Last Time: Feature Selection

- Last time we discussed feature selection:
  - Choosing set of "relevant" features.

$$X = \begin{bmatrix} & \bigcirc & \bigcirc & \\ & & & \end{bmatrix} \qquad y = \begin{bmatrix} \\ \end{bmatrix}$$

"relevant"

- Most common approach is search and score:
  - Define "score" and "search" for features with best score.
- But it's hard to define the "score" and it's hard to "search".
  - So we often use greedy methods like forward selection.
- Methods work ok on "toy" data, but are frustrating on real data…

# 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.
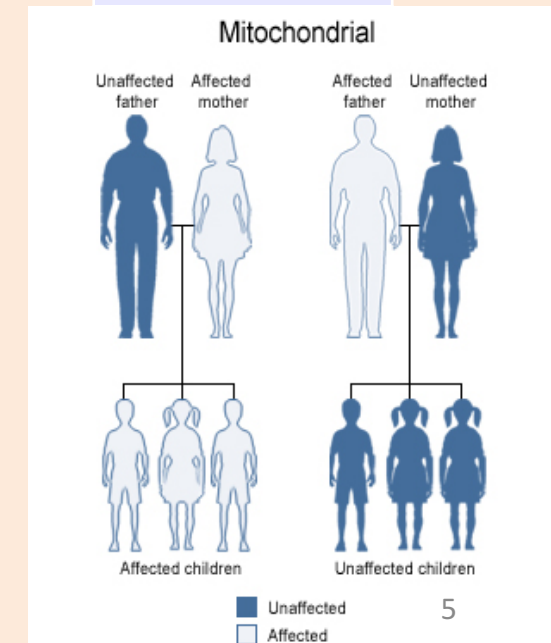    - "You almost always have the same value as your mom".

# 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.

Mitochondrial

Unaffected father | Affected mother | Affected father | Unaffected mother

Affected children | Unaffected children

■ Unaffected
□ Affected

U.S. National Library of Medicine

5

# 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 "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.
- 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 this hopeless?

- We often want to do feature selection we so have to try!

- Different methods are affected by problems in different ways.

- These "problems" don't have right answers but have wrong answers:
  - Variable dependence ("mom" and "mom2" have same information).
  - Conditional independence ("grandma" is irrelevant given "mom").

- Today we will help deal with collinearity.
- Friday we will see another powerful feature selection method.

# (pause)

# Controlling Complexity

- Usually <span style="color:red">"true" mapping from $x_i$ to $y_i$ is complex.</span>
  - Might need high-degree polynomial.
  - Might need to combine many features, and don't know "relevant" ones.
- But <span style="color:red">complex models can overfit</span>.
- So what do we do???

- Our main tools:
  - <span style="color:blue">Model averaging</span>: average over multiple models to decrease variance.
  - <span style="color:blue">Regularization</span>: add a <span style="color:green">penalty on the complexity</span> of the model.

# L2-Regularization

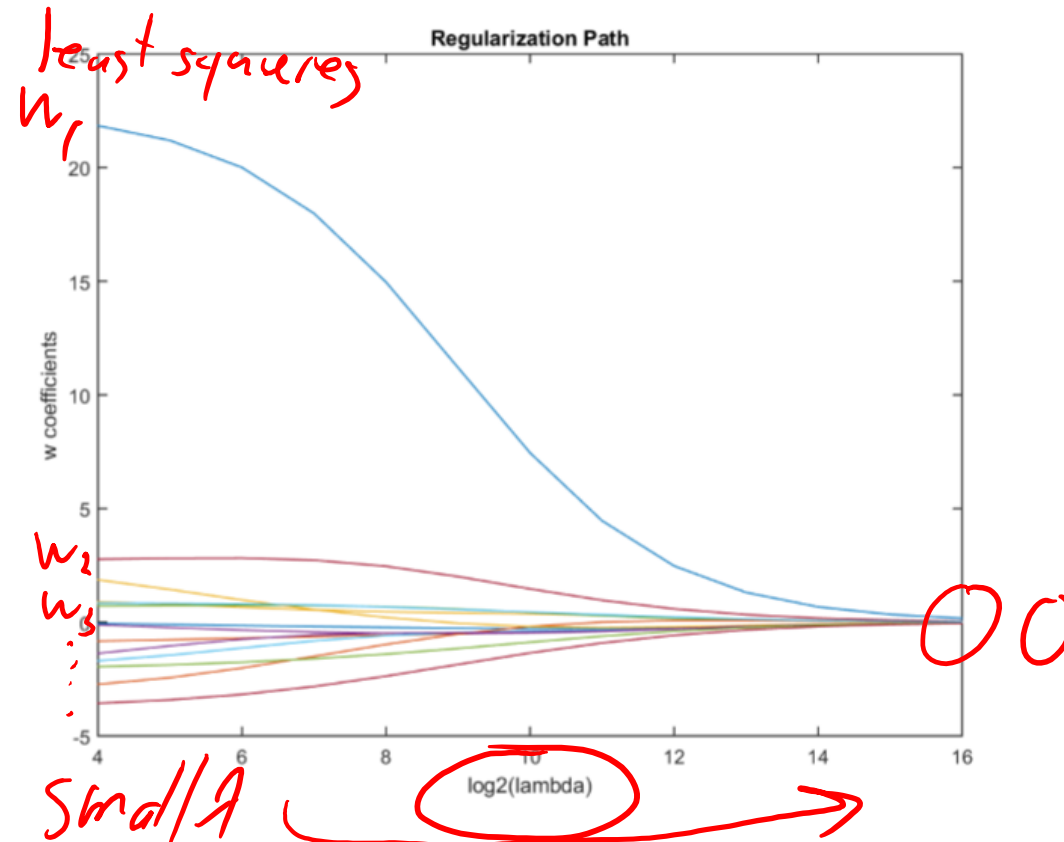- Standard regularization strategy is L2-regularization:

$$f(w) = \frac{1}{2} \sum_{i=1}^{n} (w^T x_i - y_i)^2 + \frac{\lambda}{2} \sum_{j=1}^{d} w_j^2 \quad \text{or} \quad f(w) = \frac{1}{2} \| Xw - y \|^2 + \frac{\lambda}{2} \| w \|^2$$

- Intuition: large slopes $w_j$ tend to lead to overfitting.

- So we minimize squared error plus penalty on L2-norm of 'w'.
  - This objective balances getting low error vs. having small slopes '$w_j$'.
    - "You can increase the training error if it makes 'w' much smaller."
    - Nearly-always reduces overfitting.
  - Regularization parameter $\lambda > 0$ controls "strength" of regularization.
    - Large $\lambda$ puts large penalty on slopes.

# L2-Regularization

- Standard regularization strategy is L2-regularization:

$$f(w) = \frac{1}{2} \sum_{i=1}^{n} (w^T x_i - y_i)^2 + \frac{\lambda}{2} \sum_{j=1}^{d} w_j^2 \qquad \text{or} \qquad f(w) = \frac{1}{2} \| Xw - y \|^2 + \frac{\lambda}{2} \| w \|^2$$

- In terms of fundamental trade-off:
  - Regularization increases training error.
  - Regularization decreases approximation error.

- How should you choose λ?
  - Theory: as 'n' grows λ should be in the range O(1) to $(n^{1/2})$.
  - Practice: optimize validation set or cross-validation error.
    - This almost always decreases the test error.

# Jupyter notebook demo

# Regularization Path

- Regularization path is a plot of the optimal weights '$w_j$' as '$\lambda$' varies:



- Starts with least squares with $\lambda = 0$, and $w_j$ converge to 0 as $\lambda$ grows.

# L2-regularization and the normal equations

- When using L2-regularization we can still set $\nabla$ f(w) to 0 and solve.
- Loss before:   $f(w) = ||Xw - y||_2^2$
- Loss after:   $f(w) = ||Xw - y||_2^2 + \lambda ||w||_2^2$

- Gradient before: $\nabla f(w) = X^T X w - X^T y$
- Gradient after:  $\nabla f(w) = X^T X w - X^T y + \lambda w$

- Linear system before: $X^T X w = X^T y$
- Linear system after:  $(X^T X + \lambda I)w = X^T y$
- But unlike $X^T X$, the matrix $(X^T X + \lambda I)$ is always invertible:
  - Multiply by its inverse for unique solution: $w = (X^T X + \lambda I)^{-1}(X^T y)$

# Why use L2-Regularization?

- It's a weird thing to do, but we advise "always use regularization".
  - "Almost always decreases test error" should already convince you.

- But here are more reasons:
  1. Solution 'w' is unique.
  2. $X^TX$ does not need to be invertible (no collinearity issues).
  3. Solution 'w' less sensitive to changes in X or y.
  4. Gradient descent converges faster (bigger $\lambda$ means fewer iterations).
  5. Worst case: just set $\lambda$ small and get the same performance.

# L2 regularization: not just for linear regression!

- We can use L2 (and other) regularization in a bunch of contexts
- Stay tuned for regularization later in the course:
  - Regularized linear classifiers
  - Regularized PCA
  - Regularized recommender systems
  - Regularized neural networks / deep learning
  - An interpretation of regularization in the Bayesian framework
- Overall intuition: training error isn't the goal, so why minimize it
  - Try to change the loss so that it hopefully better reflects test error

# Summary

- "Relevance" is really hard to define.
  - Different methods have different effects on what you find.
- Regularization:
  - Adding a penalty on model complexity.
- L2-regularization: penalty on L2-norm of regression weights 'w'.
  - Almost always improves test error.
  - Simple closed-form unique solution

# 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.

# 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

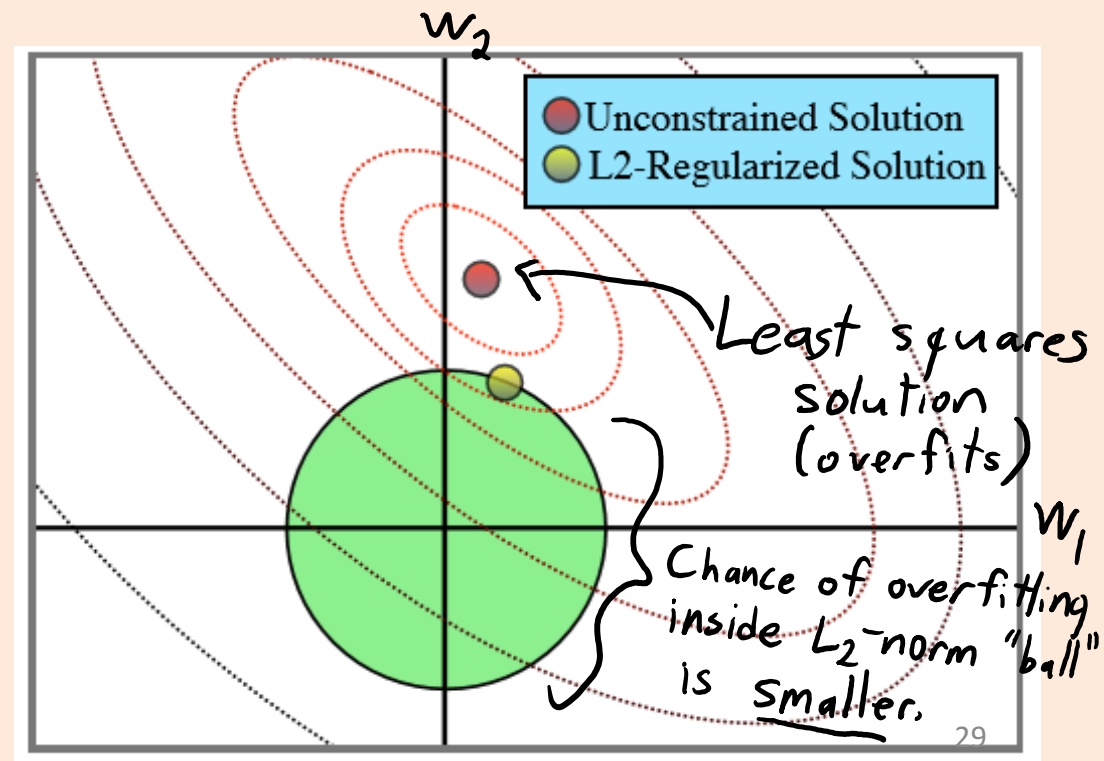| Method\Issue | Dependence | Conditional Independence | Tiny effects | Context-Specific Relevance |
|---|---|---|---|---|
| Association (e.g., measure correlation between features 'j' and 'y') | Ok (takes "mom" and "mom2") | Bad (takes "grandma", "great-grandma", etc.) | Ignores | Bad (misses features that must interact, "gender" irrelevant given "baby") |
| Regression Weight (fit least squares, take biggest $\lvert w_j \rvert$) | 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") |

# L2-Regularization

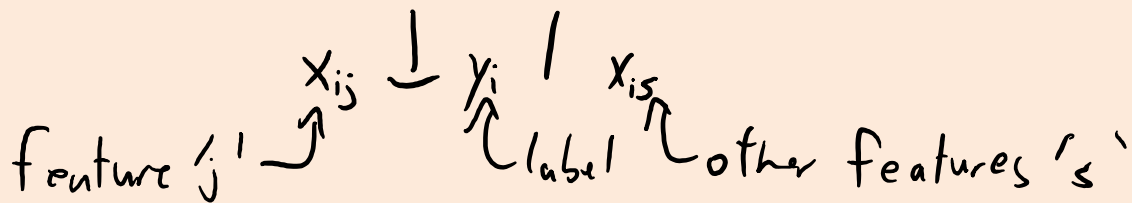- Standard regularization strategy is L2-regularization:

$$f(w) = \frac{1}{2} \sum_{i=1}^{n} (w^T x_i - y_i)^2 + \frac{\lambda}{2} \sum_{j=1}^{d} w_j^2 \quad \text{or} \quad f(w) = \frac{1}{2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2$$

- Equivalent to minimizing squared error but keeping L2-norm small.



Unconstrained Solution
L2-Regularized Solution

$w_2$

$w_1$

Least squares solution (overfits)

Chance of overfitting inside $L_2$-norm "ball" is smaller.

# 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.

$$x_{ij} \perp y_i \mid x_{is}$$

feature 'j' → $x_{ij}$   $y_i$ label   $x_{is}$ other features 's'

"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:

- Two many possible tests, "greedy" method is for each 'j' do:

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

If still dependent test $x_{ij} \perp y_i \mid x_{is}$ where 's' has one feature

If still dependent test $x_{ij} \perp y_i \mid x_{is}$ where 's' now has two features

$\vdots$

If still dependent when 's' includes all other features, declare 'j' relevant.

Often choose features to minimize dependence.

- "Association approach" is the greedy method where you only do the first test (subsequent tests remove a lot of false positives).
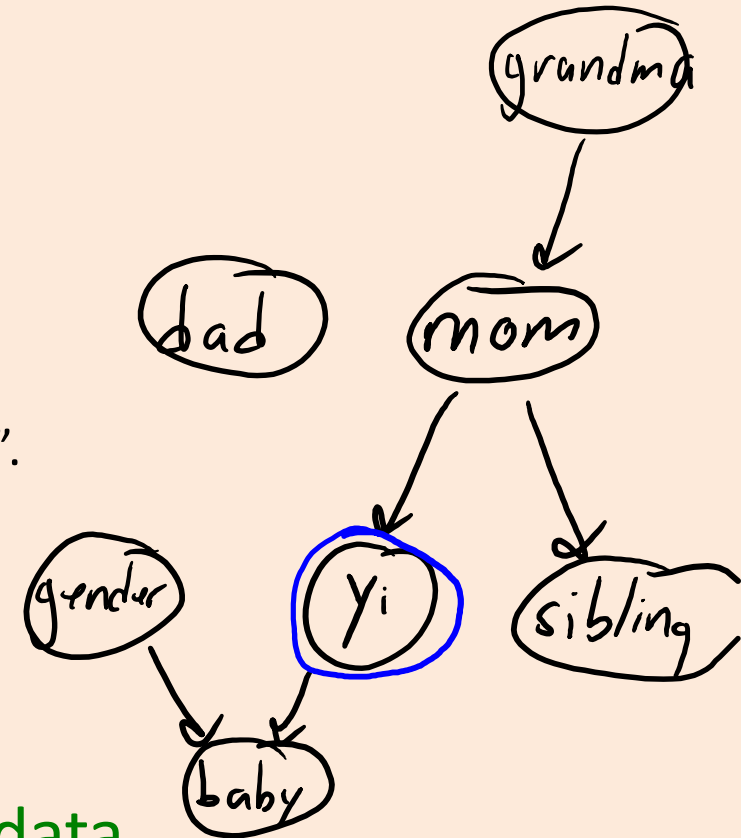
# Hypothesis-Based Feature Selection

- Advantages:
  - Deals with conditional independence.
  - Algorithm can <span style="color:green">explain why it thinks 'j' is irrelevant</span>.
  - Doesn't necessarily need linearity.

- Disadvantages:
  - Deals badly with exact dependence: doesn't select "mom" or "mom2" if both present.
  - Usual warning about <span style="color:red">testing multiple hypotheses</span>:
    - 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".
  - <span style="color:blue">This approach is sometimes better for finding relevant factors, not to select features for learning.</span>

# Causality

- None of these approaches address causality or confounding:
  - "Mom" is the only relevant 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
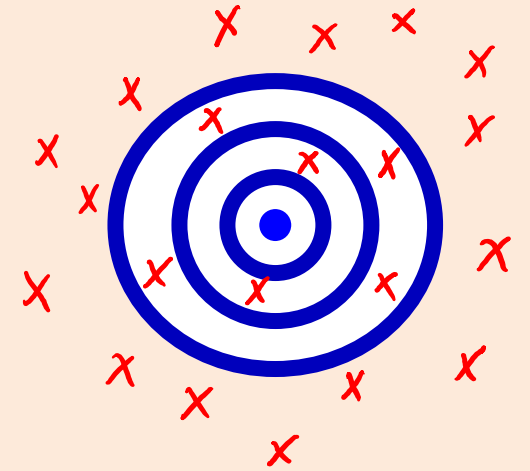
- The difference between observational and interventional data:
  - If I see that my watch says 4:45, class is almost over (observational).
  - If I set my watch to say 4:45, it doesn't help (interventional).

- 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:
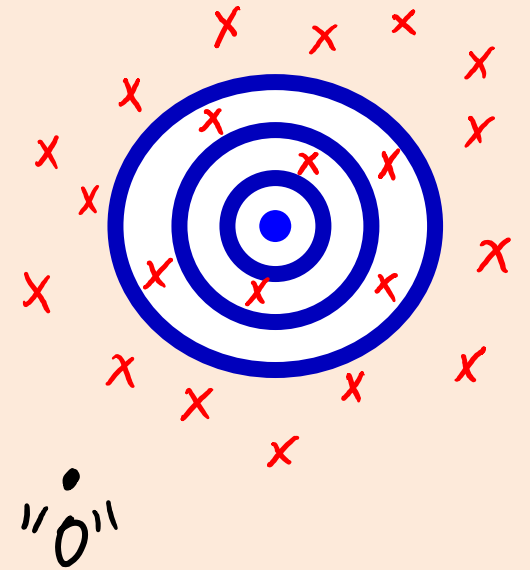  - https://www.youtube.com/watch?v=ebEkn-BiW5k

# Regularization/Shrinking Paradox

- We throw darts at a target:
  - Assume we don't always hit the exact center.
  - Assume the darts follow a symmetric pattern around center.

# Regularization/Shrinking Paradox

- We throw darts at a target:
  - Assume we don't always hit the exact center.
  - Assume the darts follow a symmetric pattern around center.

- Shrinkage of the darts :
  1. Choose some arbitrary location '0'.
  2. Measure distances from darts to '0'.

# Regularization/Shrinking Paradox

- We throw darts at a target:
  - Assume we don't always hit the exact center.
  - Assume the darts follow a symmetric pattern around center.

- Shrinkage of the darts :
  1. Choose some arbitrary location '0'.
  2. Measure distances from darts to '0'.
  3. Move misses towards '0', by *small* amount proportional to distance from 0.

- If small enough, darts will be closer to center on average.

"0"

# Regularization/Shrinking Paradox

- We throw darts at a target:
  - Assume we don't always hit the exact center.
  - Assume the darts follow a symmetric pattern around center.

- Shrinkage of the darts :
  1. Choose some arbitrary location '0'.
  2. Measure distances from darts to '0'.
  3. Move misses towards '0', by *small* amount proportional to distance from 0.

- If small enough, darts will be closer to center on average.

Visualization of the related higher-dimensional paradox that the mean of data coming from a Gaussian is not the best estimate of the mean of the Gaussian in 3-dimensions or higher: https://www.naftaliharris.com/blog/steinviz