Thirteen ways of looking at a default

MAS.622J Final Project - Prosper.com

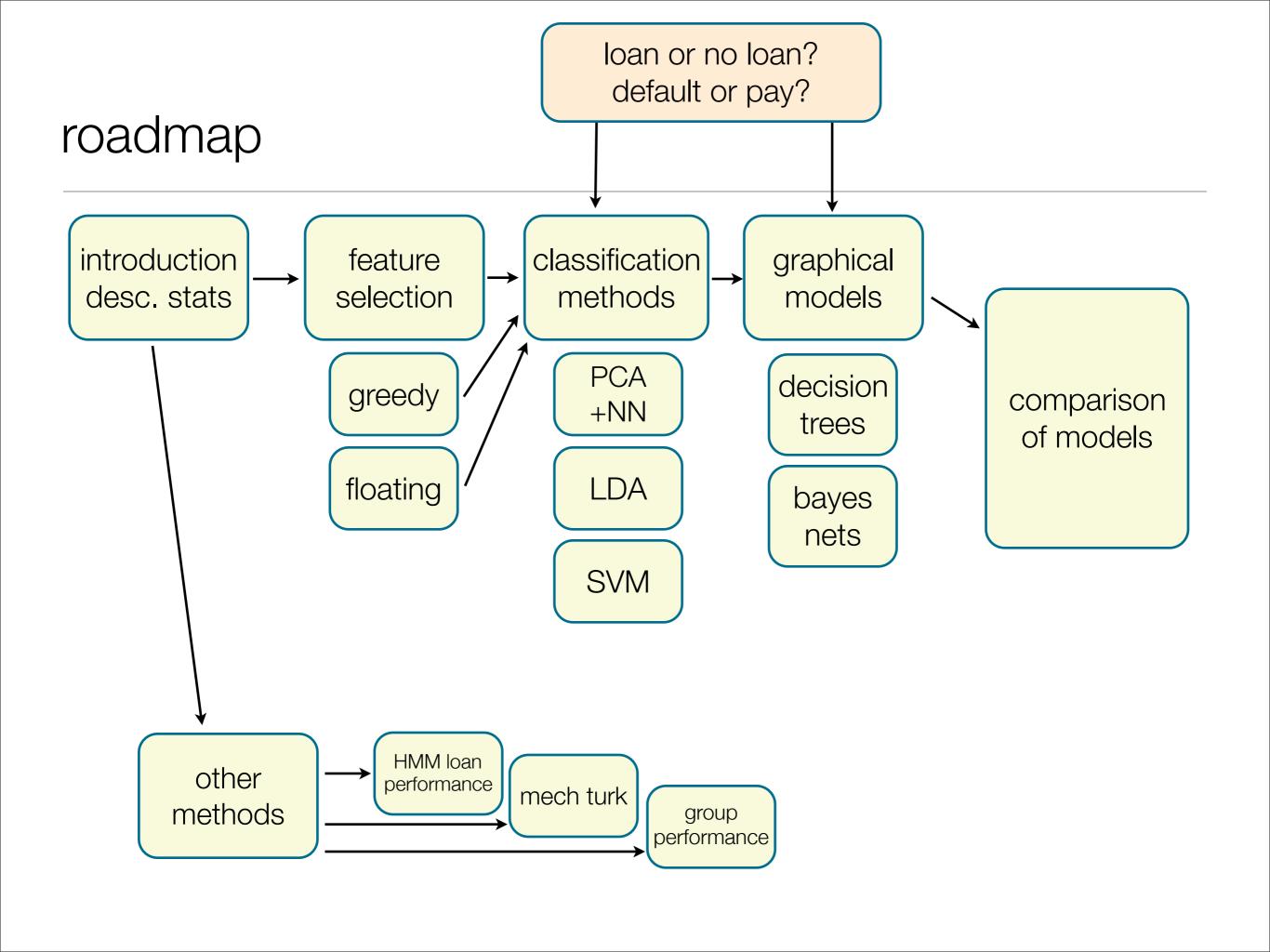
Charlie deTar
Coco Krumme
Ernesto Martinez-Villalpando
Matt Aldrich

review of questions

- What is the likelihood of a listing becoming a loan?
- What is the likelihood of a loan being paid on time?
- Which features best predict these outcomes?

why this is important

- Borrowers want to maximize likelihood of getting a loan, minimize interest rate
- Lenders want to invest in loans that maximize returns, minimize probability of default or late payment
- Prosper.com wants to maximize revenues by increasing loan conversion, decreasing default rate
- The research questions are deep: is peer-to-peer lending a viable model? How much do social factors matter? How do lenders make decisions? What models best capture loan dynamics? Can peer-to-peer be modeled with the same precision as bank loans? How can human classification aid machine learning algorithms? etc etc



Charlie De Niro Matt Stiller Ernesto Hoffman Coco Streissand





UNIVERSAL PICTURES AND DREAMWORKS PICTURES PRESENT A TRIBECA / EVERYMAN PICTURES PRODUCTION A JAY FOACH FUM
ROBERT DENIED BEN STILLER DUSTIN HOFFMAN AND BARBRA STRESAND WEET THE FOCKER'S BYTHE DANNER TEST POLD WERTANDY NEWMAN BERNEGARD RAMASY AND POL

THE JON POLL LEE HAVALL "MUSTUREDUSTY SMITH AUGUSTA JOHN CHWARTZHAN ASS, AUGUSTA NANNY TENENBAUM ALMY SAVES HAVE DEMANDED FOR CLEHNA & MARY BUTH CLARK

"WILL JON POLL LEE HAVALL "MUSTUREDUSTY SMITH AUGUSTA JOHN FOR CHWARTZHAN ASS, AUGUSTA NANNY TENENBAUM ALMY SAVES HAVE DEMANDED FOR CHIEF JOHN FOR CHIEF JOHN

Prosper data tables

Bid

- Amount
- Minimum Rate
- Listing Status
- ...

Loan

• Kev

• Name

Credit Grace

Category

Hierarchy

- Borrow Rate
- Debt to Income Ratio
- ٠...

Group

- Member Key
- Group Rating
- City
- ٠...

Listing

- Amount Funded
- Amount Remaining
- Bid Count
- ...

Credit Profile

- Amount Delinquent
- Bankcard Utilization
- Credit Grade
- ٠...

Member

- Key
- Friend Member Keys
- Group Key
- ٠ ...

Loan Performance

- Cycle Number
- DPD (Date Past Due)
- Net Defaults
- ٠..

Marketplace

- Groups Count to Date
- Interest Rates Table
- Loans Closed Count
- ٠ ...

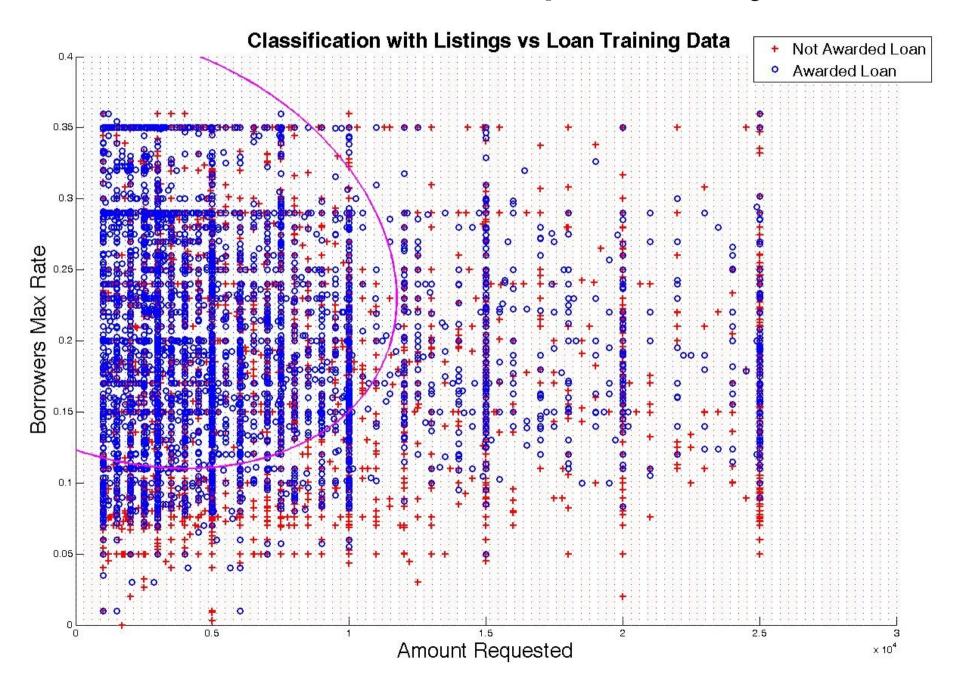
"the original 11"

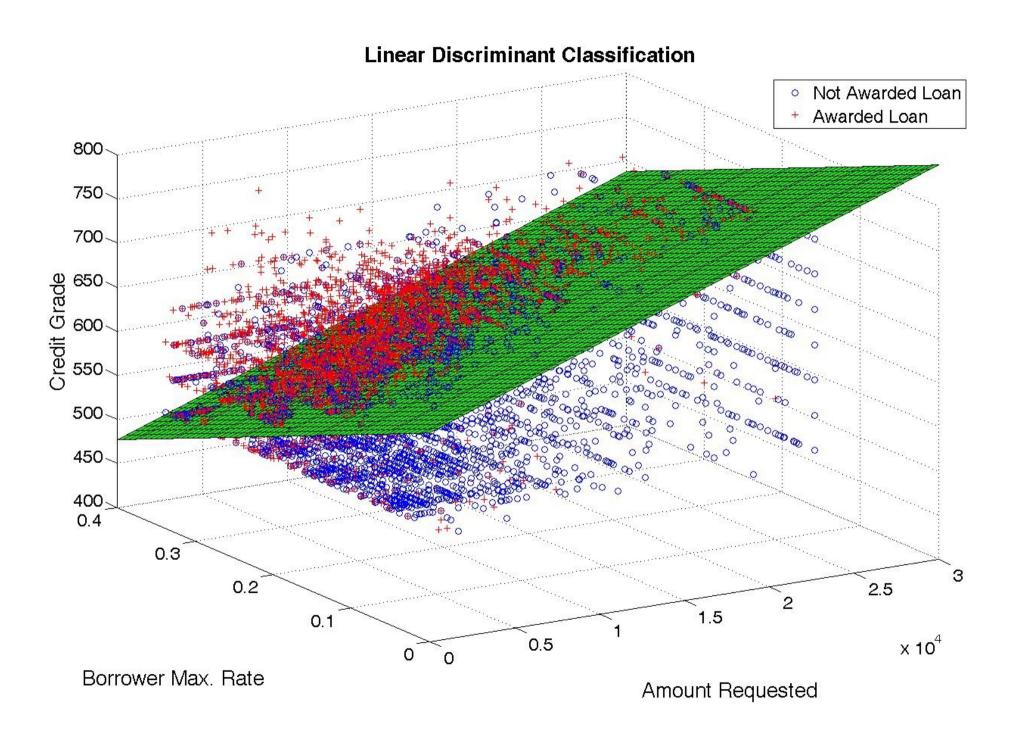
- Amount Requested
- Bid Count *
- Borrower Rate
- Credit Grade
- Debt to Income Ratio
- Group Key

- Has an image
- Current delinquencies
- Delinquencies last 7 years
- Open credit lines
- Income

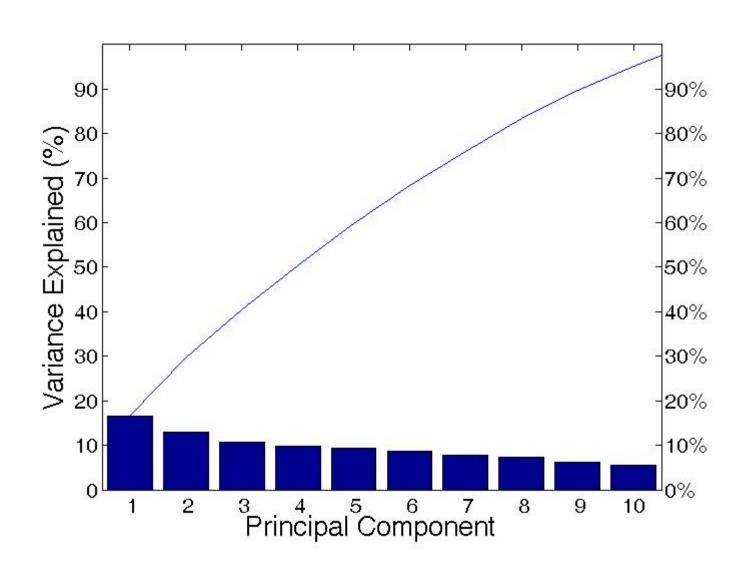
^{*} Not used in Loans vs. unfunded listings classification

PCA shows separability





Variance and principle components



Diverse... and non-numerical

- Textual fields
- Number... or null
- Binary next to thousands



Hi and thank you for looking at my post. I currently own a small 3 employee business in Minnesota, I started the business about 3 years ago and it is really taking off. We currently have over \$100,000 in inventory and are looking to hire more employees. I would like this loan to actually buy even more inventory and also just to have fun with Prosper and use it, I love lending people money on here, it is way more fun then the stock market. Our current sales are about \$360,000/yr. but we are on track for \$500,000 this next fiscal year.

A little info about myself: I am married to a wonderful woman and we have a baby on the way. We have a beautiful big brown Newfoundland named Tank and a Persian cat named Goo who hates me. My hobbies include playing hockey, flying small private planes, and building stuff around my house.

Thank you!

Descriptions

Loans	Listings	Difference	Words
0.44	0.58	0.14	cards and other
0.29	0.42	0.14	monthly expenses housing
0.44	0.58	0.14	clothing household expenses
0.41	0.54	0.13	and other loans
0.40	0.50	0.44	
0.42	0.56	0.14	other expenses
0.29	0.43	0.14	expenses housing
0.44	0.57	0.13	clothing household
0.45	0.58	0.13	car expenses

Titles

Loans	Listings	Difference	Words
0.024	0.014	0.010	high interest credit
0.022	0.015	0.007	off high interest
0.021	0.016	0.005	credit card debt
0.009	0.004	0.005	interest credit card
0.062	0.04	0.022	credit card
0.049	0.032	0.017	high interest credit
0.015	0.003	0.012	in prosper
0.026	0.015	0.011	interest card

"the 96"

Amount delinquent null? Amount delinquent Bankcard utilization null? Bankcard utilization Current credit lines null? Current credit lines Current delinquencies null? Current delinquencies Delinquencies last 7 years null? Delinquencies last 7 years Income Length status months Open credit lines null? Open credit lines Revolving credit balance null? Revolving credit balance Amount requested Borrower maximum rate Listing category Credit grade Debt to income ratio null? Debt to income ratio "prosper" in description? "as" in description? "clothing" in description? "household" in description? "housing" in description? "card" in description? "entertainment" in description? "is" in description? "with" in description?

"an" in description? "other expenses" in description? "clothing household" in description? "car expenses" in description? "household expenses" in description? "other loans" in description? "phone cable" in description? "and other" in description? "food entertainment" in description? "cards and" in description? "cards and other" in description? "monthly expenses housing" in description? "clothing household expenses" in description? "and other loans" in description? "phone cable internet" in description? "credit cards and" in description? "monthly net income" in description? "a good candidate" in description? "good candidate for" in description? "my financial situation" in description? Funding option Is borrower homeowner null? Is borrower homeowner? "bills" in title "prosper" in title "card" in title "interest" in title "credit" in title? "credit card" in title? "high interest" in title?

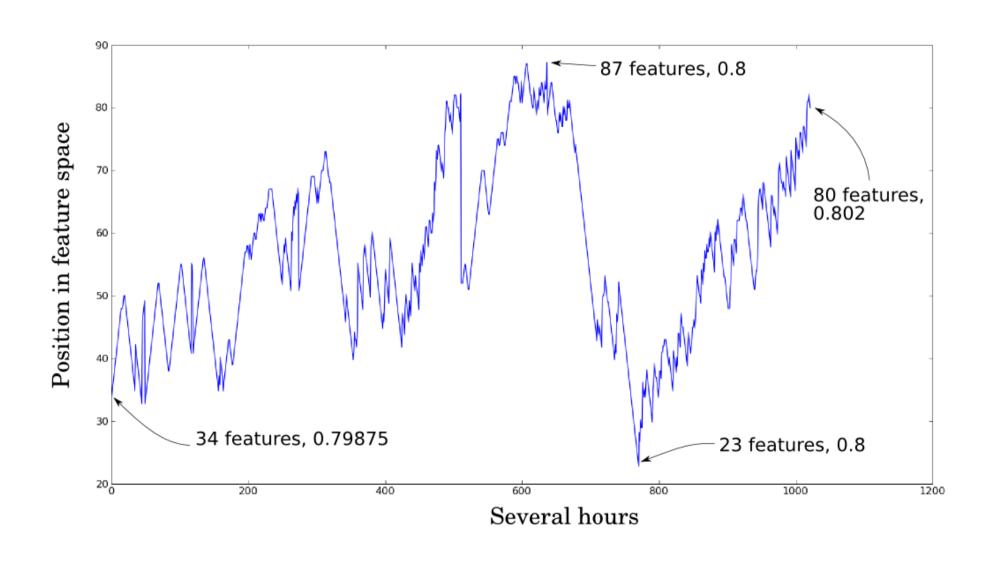
"in prosper" in title?

"interest credit" in title? "pay off" in title? "high interest credit" in title? "off high interest" in title? "credit card debt" in title? "interest credit card" in title? "off credit cards" in title? # of non-alphanumeric characters in title number of CAPS in title Member endorsements null? "i" in member endorsements? "and" in member endorsements? "a" in member endorsements? "to" in member endorsements? "the" in member endorsements? "i have" in member endorsements? "is a" in member endorsements? "this loan" in member endorsements? "will be" in member endorsements? "he is" in member endorsements? "i have known" in member endorsements? "he is a" in member endorsements? "this is a" in member endorsements? "i will be" in member endorsements? "i've known" in member endorsements? Number of "friends" Is member of a group? Has an image in description? Is a borrower Is a Group Leader Is a lender Is an institutional lender

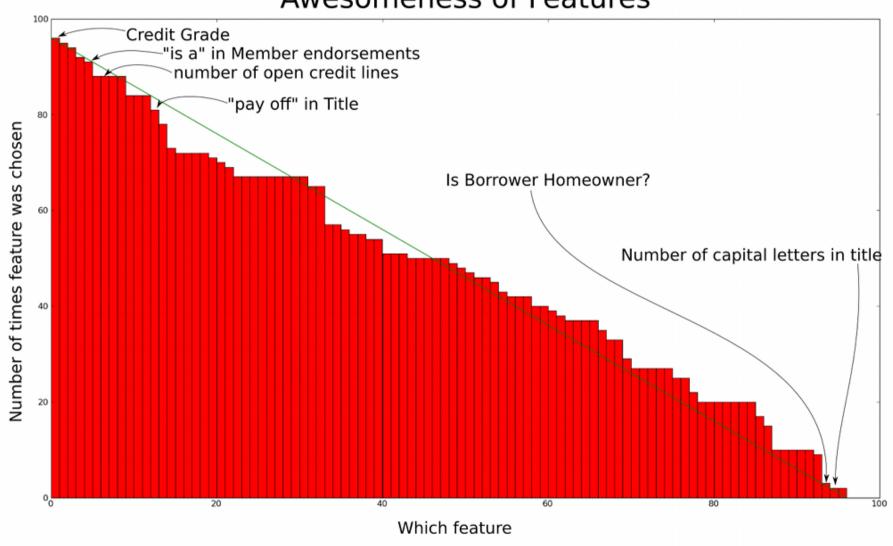
Floating Feature Search

- Linear Discriminant as evaluation function
- Lots of samples, lots of features:
 - 96 features
 - 300,000 listings
- sssssslllllllooooooowwwwwwww. Must:
 - limit the number of features hence forwards floating search. Decreases optimality.
 - decrease the number of samples (increases bias)

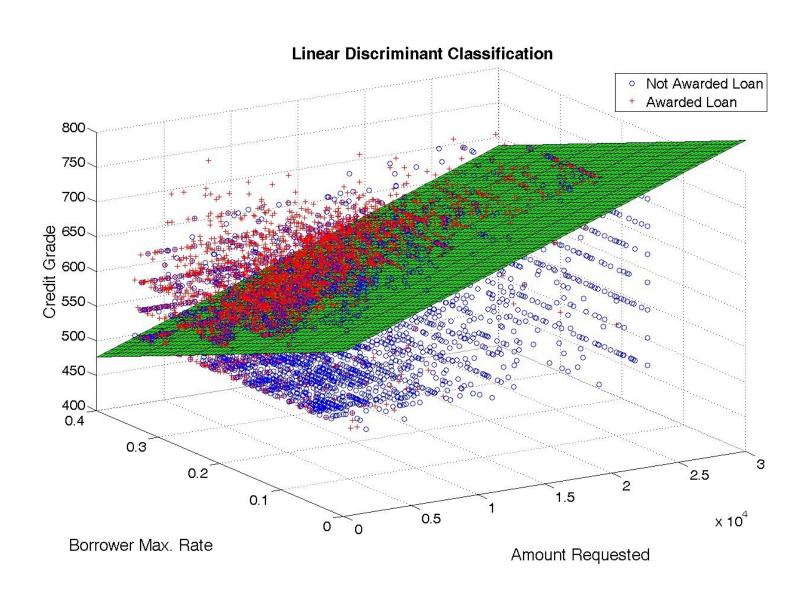
Feature Crawl



Awesomeness of Features



Classification Performance

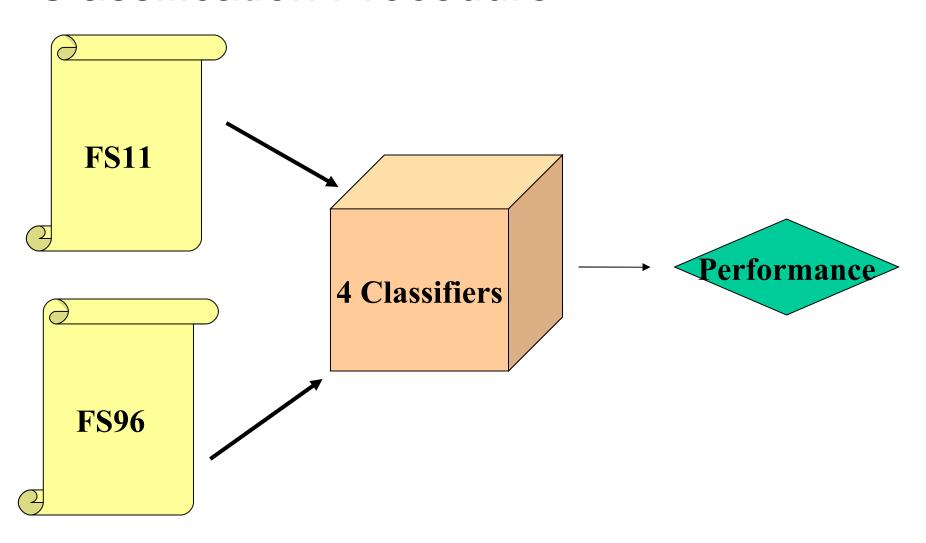


Classification Overview

- Review of Methods
- Discussion of prior probability, implications when viewing results
- Summary of Results, Tables
- Classification Improvements
- Tying it Back to P2P Lending
- Detailed performance data is given on project website.

What We Did

Classification Procedure



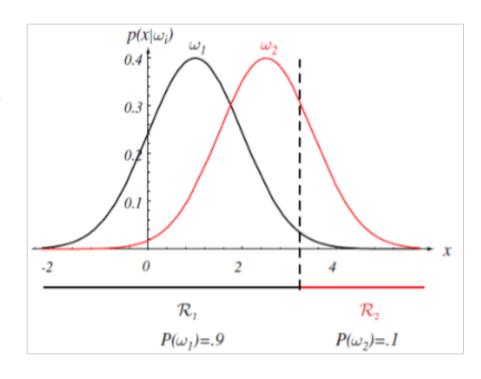
Review of Methods

- LDA pseudoinverse (mse)
- PCA+LDA reduce dims, classify using some number of principal comp. (70%/30% 10F CV)
- SVM map higher to higher dim space
 - Linear, 1 norm (smo) soft margin (slack) (70%/30%, no CV)
- Neural Networks high degree of freedom (hidden nodes)
 - Feedforward, Tr:conjugate gradient descent. 1 HL {5,10,20 hidden nodes} 70%tr,15%val, 10%test
- Libraries: Matlab SVM, NN Tools

Listing / Loan Discussion

- Prior probability leads to classifiers that favor one class
- In comparing classifiers using stratified sampling FN rate is large contributor of error
- Results are given for both stratified sampling & equal priors
- PCA+LDA classification attempts to separate these distributions in a lower dim

TP	FN
FP	TN



From: Richard O. Duda, Peter E.

Hart, and David G. Stork, Pattern Classification. Copyright c 2001 by John Wiley & Sons, Inc.

Brief Summary of Results

- Listing\Loan (Stratified) best performing classifier
 - FS10, Neural Network (10 hidden nodes)
 - 14% error (prevalence of C1: 16.8%)
 - 2000 samples, 20 Hidden, 75%Tr, 15%Val, 10%Test
 - Different prior probabilities of FS10 & FS96
 - Effect is that error is mostly FN approx. = prevalence of c1
- Listing\Loan (Equal) best performing classifier
 - FS10, FS96 Neural Network (20 hidden nodes)
 - 18%, 16% respectively
 - 2000 samples, 20 Hidden, 75%Tr, 15%Val, 10%Test
- Default\No Default
 - FS11, FS96 Neural Network (20 hidden nodes)
 - 26%,15% respectively
 - 2000 samples, 20 Hidden, 75%Tr, 15%Val, 10%Test

Lessons \ Future Improvements

- Neural Nets were a good match, surprising?
 - Not really, given a number of hidden nodes (degrees of freedom), arbitrarily complex decision boundaries can be found – great for high dimension feature vectors.
- Effect of adding additional features
 - For any method of classification, data suggests additional dimensions improve accuracy, complexity not worth the effort. We are talking about ~4% less error.
- Satisficing LDA+PCA good enough
 - No matter what method was used results ~84-79% correct
- Real world data != equal priors
 - Feature search should seek to minimize FP,FN better separation, more realistic classifiers for Loan/No Loan

Take Home Message

- What is all this really good for anyways?
 - Designed several classifiers performing > 80% accuracy (that's great but...)
 - Goal is not to make the world's best performing classifier, rather the data can be classified. (clustering-> classification -> intuitive models)
- A Strategy for Borrowers?
 - We demonstrated that there *are* features that can separate the data, what is your strategy to improve chances...
 - Classification results not quite satisfying and tractable
 - Coco & Ernesto build on these results demonstrate models that make intuitive sense
- A Strategy for Lenders an Investment Tool?
 - We demonstrated to a reasonable accuracy that features indicative of defaulting exist. Is this better than intuition? (a machine classifier doesn't say you should get a loan because you got divorced, experienced a disaster, etc)
 - Problems with this strategy Someone of dubious repayment potential gets a loan. She then repays because she won the lottery. Outside events not taken into account.
 - A temporal analysis to examine loan performance is required. "Now that you have a loan lets see what you do." Stay tuned for HMMs

Loan / No Loan Summary

Method	FS96 E	FS10 E	#FS9 6	#FS1 0	Ratio FS96	Ratio FS10	Tr/(Val)/ Te	CV FS96	CV FS11
SVM (All)	0.085	0.157	2000	2000	0.085	0.168	70/30	None	None
SVM TOP 10	0.098	null	2000	2000	0.085	0.168	70/30	None	None
PCA+LDA (10 P.C.)	0.259	0.223	5000	5000	0.085	0.168	70/30	10 Fold CV	10 Fold CV
NN ALL (20 Hidden)	0.085	0.143	2000	2000	0.085	0.168	75/15/10	N/A	N/A
NN ALL (10 Hidden)	0.085	0.141	2000	2000	0.085	0.168	75/15/10	N/A	N/A
NN ALL (5 Hidden)	0.085	0.145	2000	2000	0.085	0.168	75/15/10	N/A	N/A
LDA 30	0.085	null	10000	null	0.085	null	70/30	None	null
LDA 10	0.085	0.141	10000	10000	0.085	0.168	70/30	None	10 Fold CV
LDA 5	0.085	null	10000	null	0.085	null	70/30	None	null

Note: while the error
In this case is high: the
FN classification is better
due to pca dim reduction

Stratified sampling

Method	FS96 E	FS10 E	#FS9 6	#FS1 0	Ratio FS96	Ratio FS10	Tr/(Val)/ Te	CV FS96	CV FS11
SVM (All)	0.192	0.230	2000	2000	0.500	0.500	70/30	None	None
SVM TOP 10	0.212	null	2000	2000	0.500	0.500	70/30	None	None
PCA+LDA (10 P.C.)	0.259	0.223	5000	5000	0.500	0.500	70/30	10 Fold CV	10 Fold CV
NN ALL (20 Hidden)	0.160	0.185	2000	2000	0.500	0.500	75/15/10	N/A	N/A
NN ALL (10 Hidden)	0.183	0.191	2000	2000	0.500	0.500	75/15/10	N/A	N/A
NN ALL (5 Hidden)	0.212	0.200	2000	2000	0.500	0.500	75/15/10	N/A	N/A
LDA 30	0.213	null	10000	null	0.500	null	70/30	None	null
LDA 10	0.259	0.220	10000	10000	0.500	0.500	70/30	None	10 Fold CV
LDA 5	0.224	null	10000	null	0.500	null	70/30	None	null

Equal Sampling

Default / No Default Summary

Method	FS96 E	FS11 E	#FS9 6	#FS1 1	Ratio FS96	Ratio FS11	Tr/ (Val)/Te	CV FS96	CV FS11
SVM (All)	0.190	0.270	2000	2000	0.500	0.500	70/30	None	None
SVM TOP 10	0.210	null	2000	2000	0.500	0.500	70/30	None	None
PCA+LDA (10 P.C.)	0.197	0.250	2000	2000	0.500	0.500	70/30	10 Fold CV	10 Fold CV
NN ALL (20 Hidden)	0.152	0.264	2000	2000	0.500	0.500	75/15/10	N/A	N/A
NN ALL (10 Hidden)	0.156	0.271	2000	2000	0.500	0.500	75/15/10	N/A	N/A
NN ALL (5 Hidden)	0.164	0.273	2000	2000	0.500	0.500	75/15/10	N/A	N/A
LDA 30	0.273	null	1000	null	0.500	null	70/30	None	null
LDA 11	0.257	0.240	1000	2000	0.050	0.500	70/30	None	10 Fold CV
LDA 5	0.248	null	1000	null	null	null	70/30	None	null

Bayesian Network

- Nine Quantized Features from Floating Selection Set:
 - A. Amount Delinquent (Low, High)
 - B. Open Credit Lines (Low, Med, High)
 - C. Amount Requested (Low, Moderate, High, Very High)
 - D. Borrower's Max Rate (Low, Moderate, High, Very High)
 - E. Credit Grade (Poor, Average, Good, Very good)
 - F. Debt to Income Ratio (Low, Med, High)
 - G. 'Good Candidate' (True, False)
 - H. Funding Option (True, False)
 - I. Endorsement (True, False)

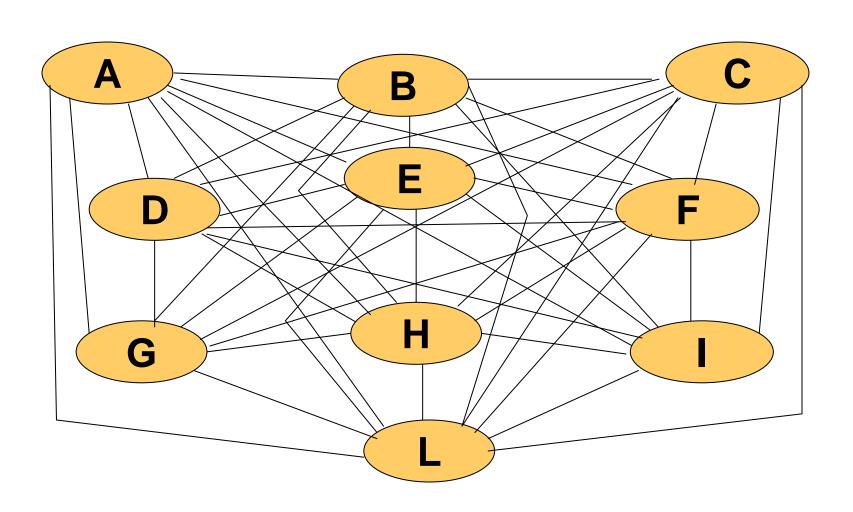
Structure Learning

- Methods:
 - Exhaustive Search: PC Algorithm
 - Score-Based:MCMC; K2, Greedy Search
- Challenges:
 - 10 Nodes = 4.2 x 10^18 Directed Acyclic Graphs!!!!!!!

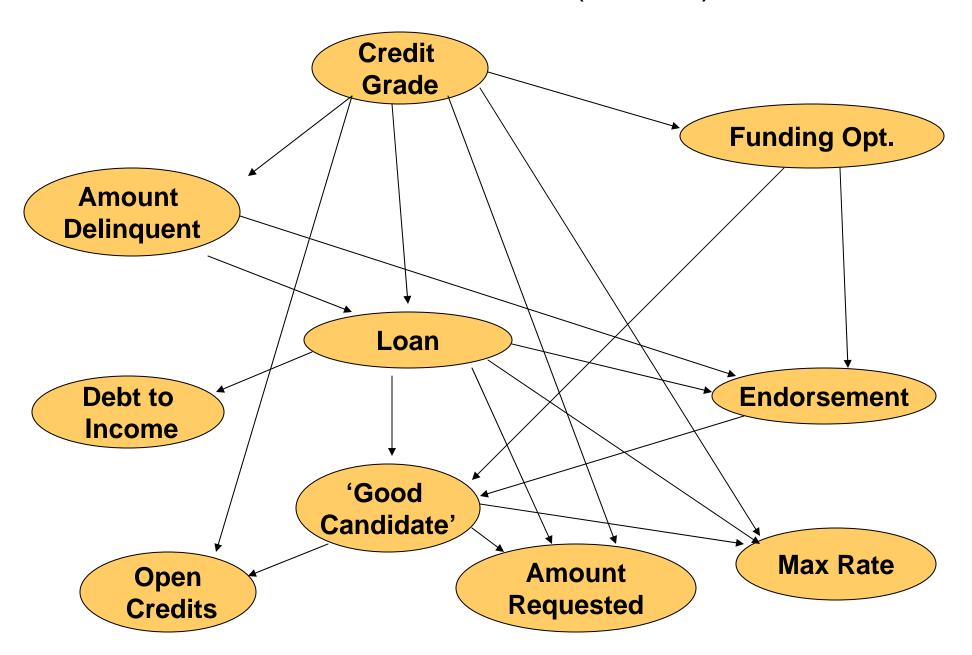
$$r(n) = \sum_{i=1}^{n} (-1)^{i+1} \binom{n}{i} 2^{i(n-i)} r(n-i) = n^{2^{\mathcal{O}(n)}}$$

- PC algorithm... Overflow!!!
- Overfitting??

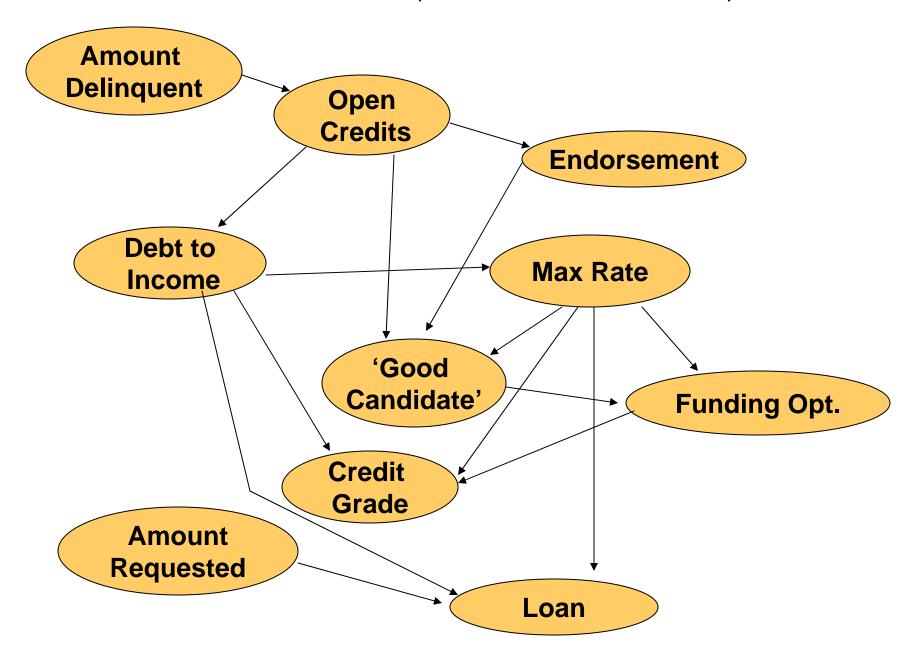
Complete Graph



Learned Structure (MCMC)



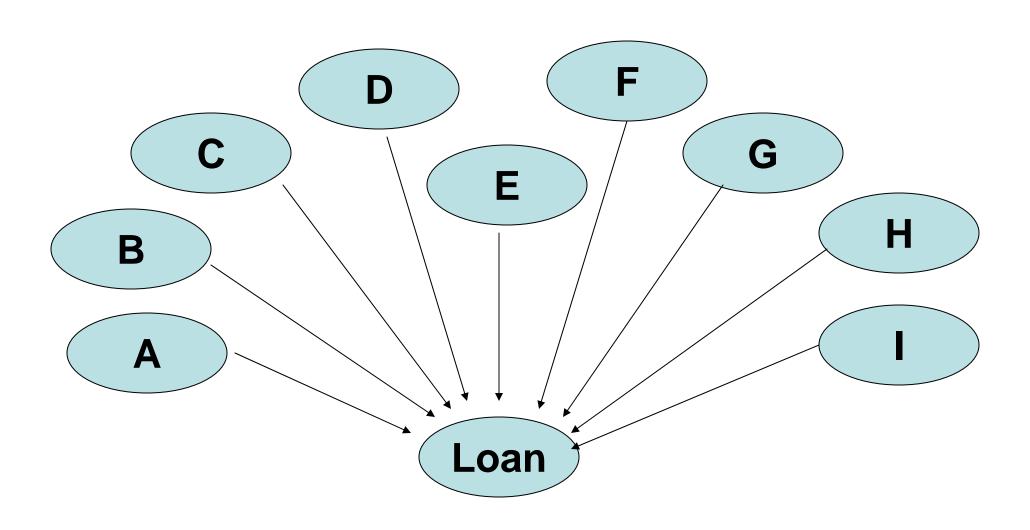
Learned Structure (Chow-Liu Trees & K2)



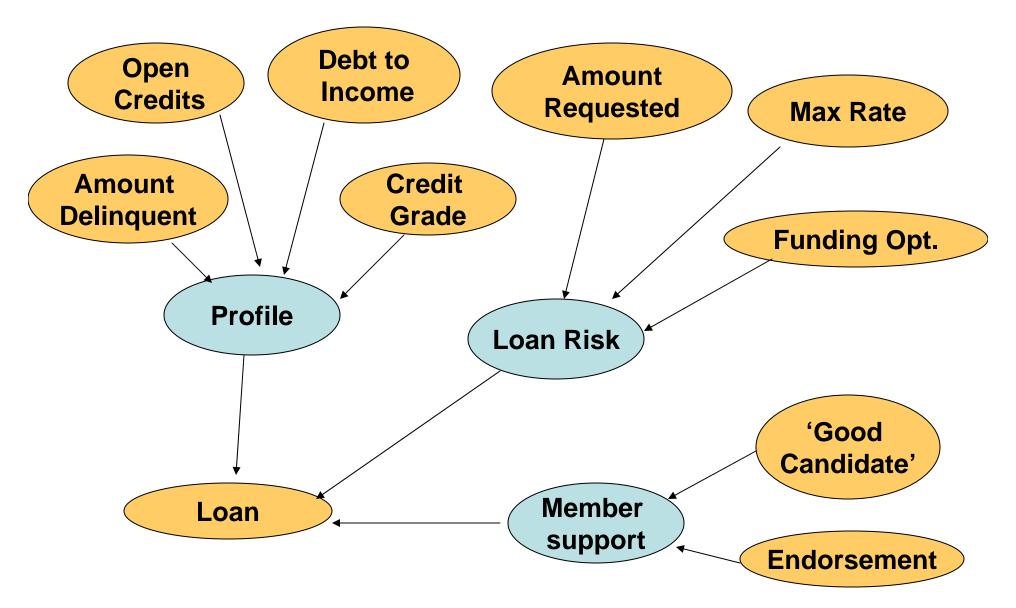
Building Other Models

- Models:
 - Naïve Bayes Classifier
 - Belief Structures
 - Noisy Functional Dependence Models
- Parameter Estimation (complete data set)
 - Batch Learning: MLE & Bayesian Estimation (Maximum a posteriori parameters)
 - MAP decision rule (classification)

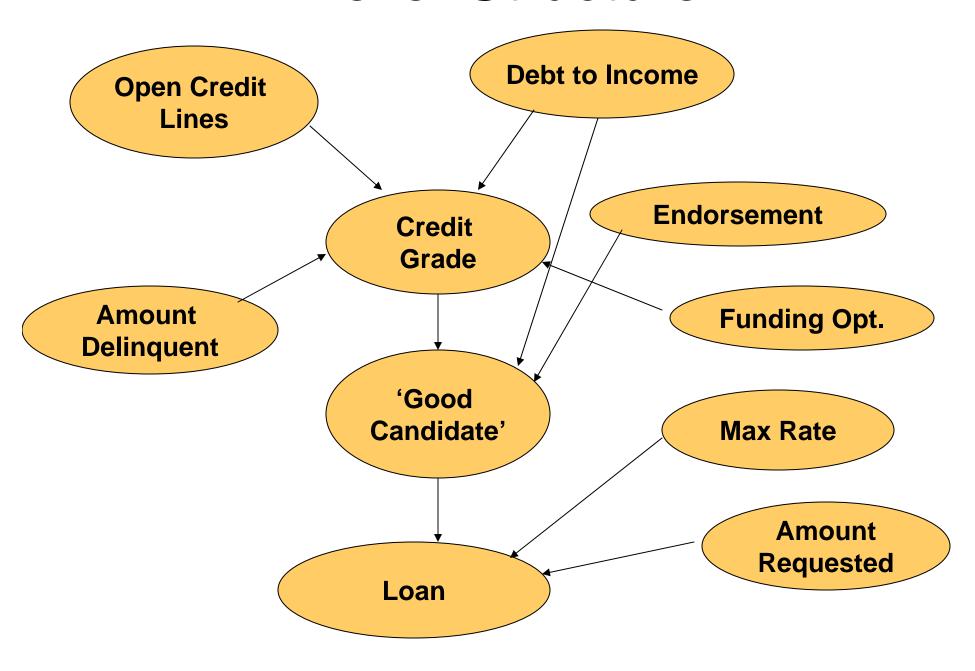
Naïve Bayes Classification



Noisy Functional Dependence

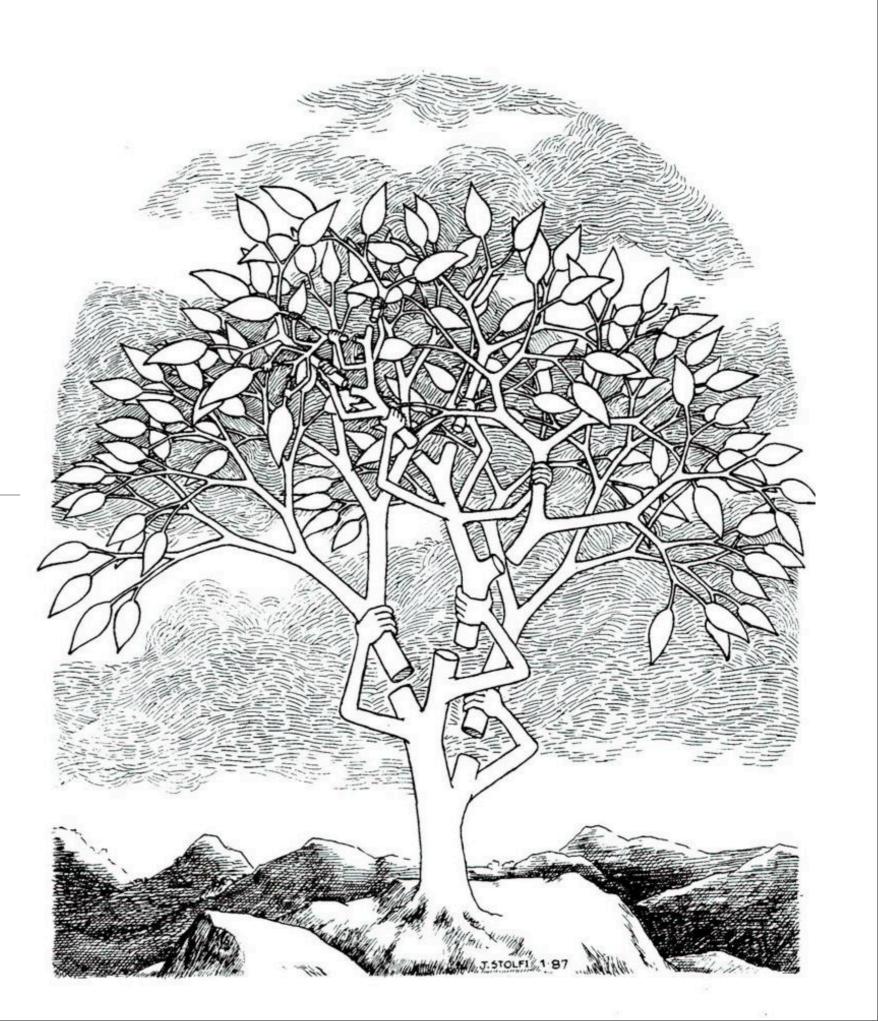


Belief Structure



Bayes Net Model	BIC score (x 10^5)	Clsf. Performance
Learned Structure (MCMC)	-1.33	0.76
Learned Structure (Chow-Liu & K2)	-1.34	0.7525
Naïve Bayes Net	-1.83	0.7580
Belief structure	-1.42	0.5620

Decision Trees



Decision Trees: questions

- BORROWER: will my loan get funded?
 - (how much should I borrow? what interest rate should I set?)
- LENDER: if I fund this loan, will I be paid back?
 - (what features best predict default? which loans should I fund?)

Decision Trees: methodology

- FEATURES: experimented with various sets
 - Greedy 11 (eliminated #bids)
 - sets of 2 4 6 8 10 features
- NODE SPLIT THRESHOLD: 2 6400

probability i belongs to class j

minimize Gini impurity

$$I_G(i) = 1 - \sum_{j=1}^m f(i,j)^2 = \sum_{j \neq k} f(i,j)f(i,k)$$

to zero when all samples part of single target category

• PRIORS: tried with / without prior probabilities [13% loan, 87% no loan]

Decision Tree: analysis

Variables:

- FEATURE SET
- NODE SPLIT threshold

Tests:

- SENSITIVITY: reserved 10%, 10 iterations
- ERROR RATES: total, FP, FN



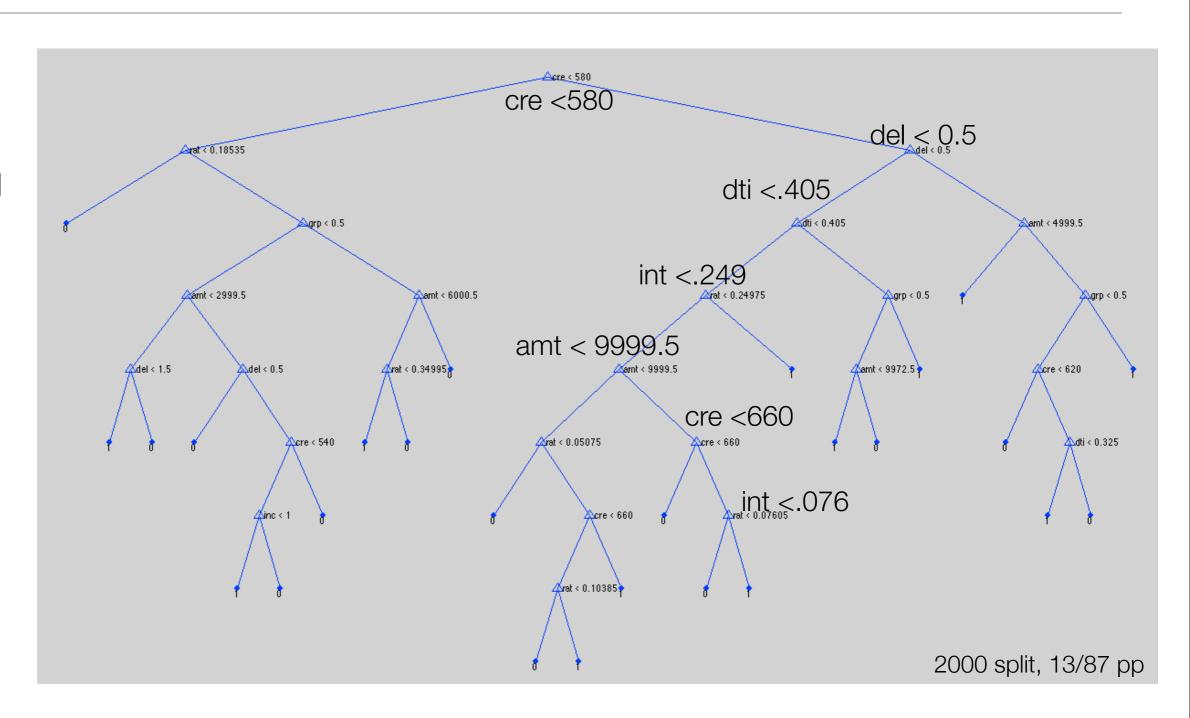
Pick best feature sets & split threshold to:

minimize variance across iterations minimize error

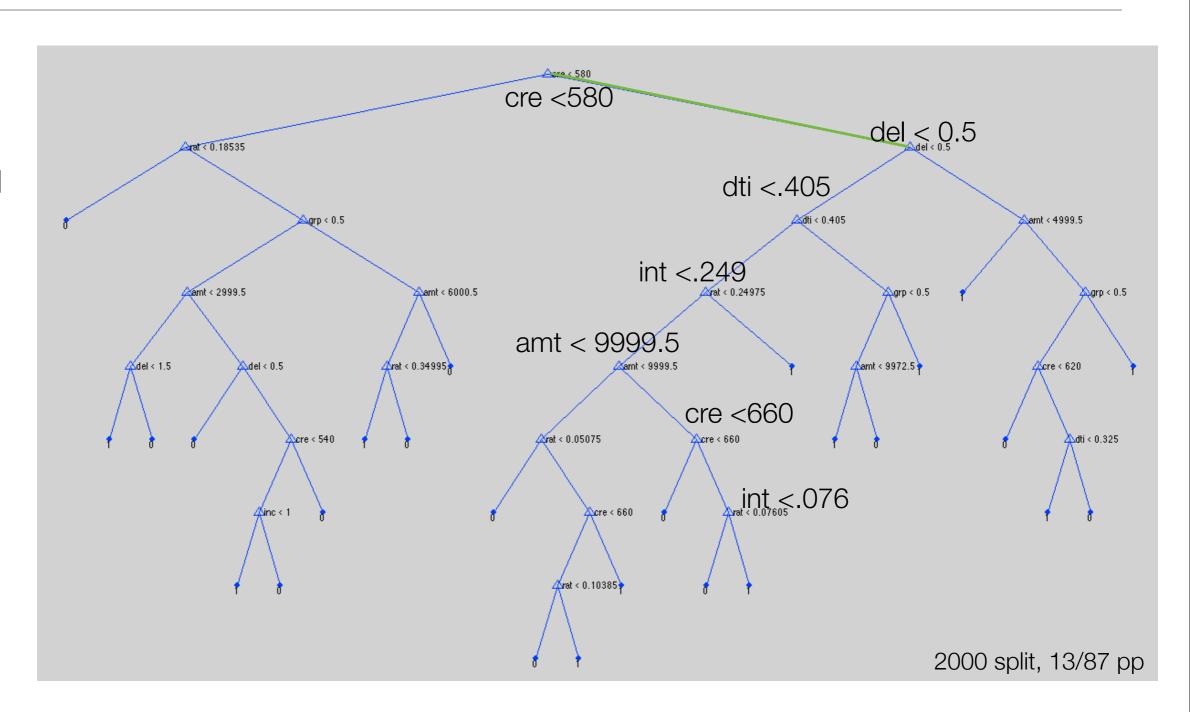
now entering: TINY FONT TERRITORY



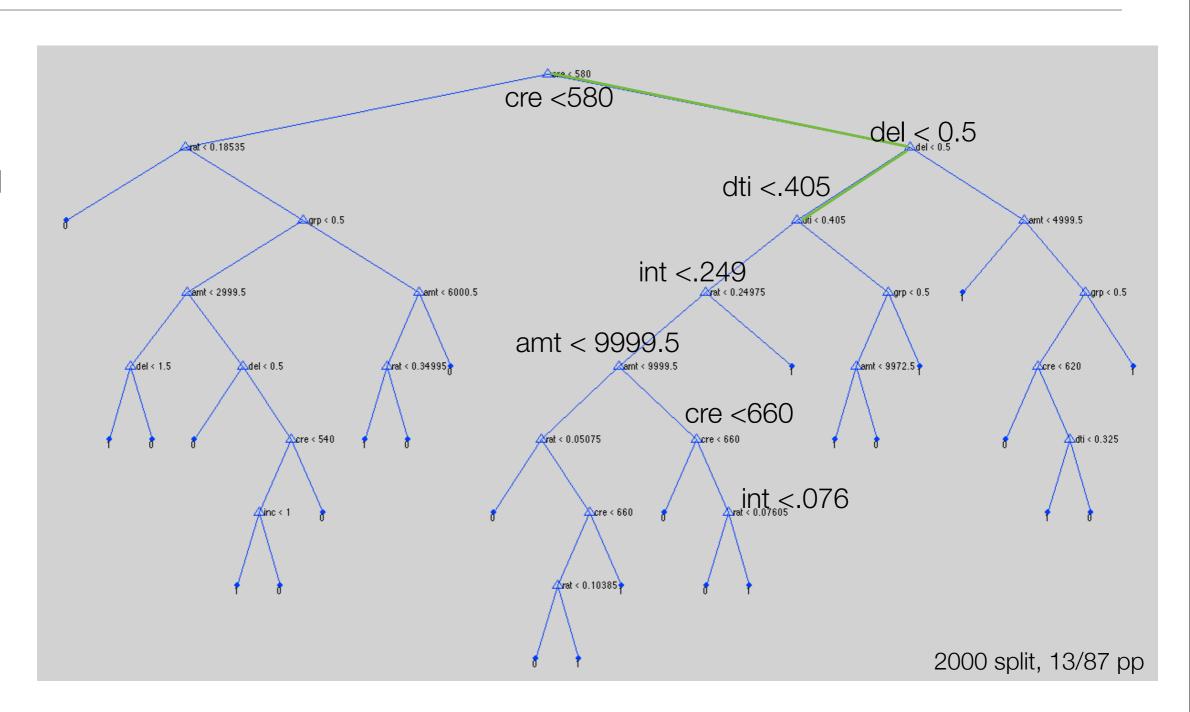
- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan



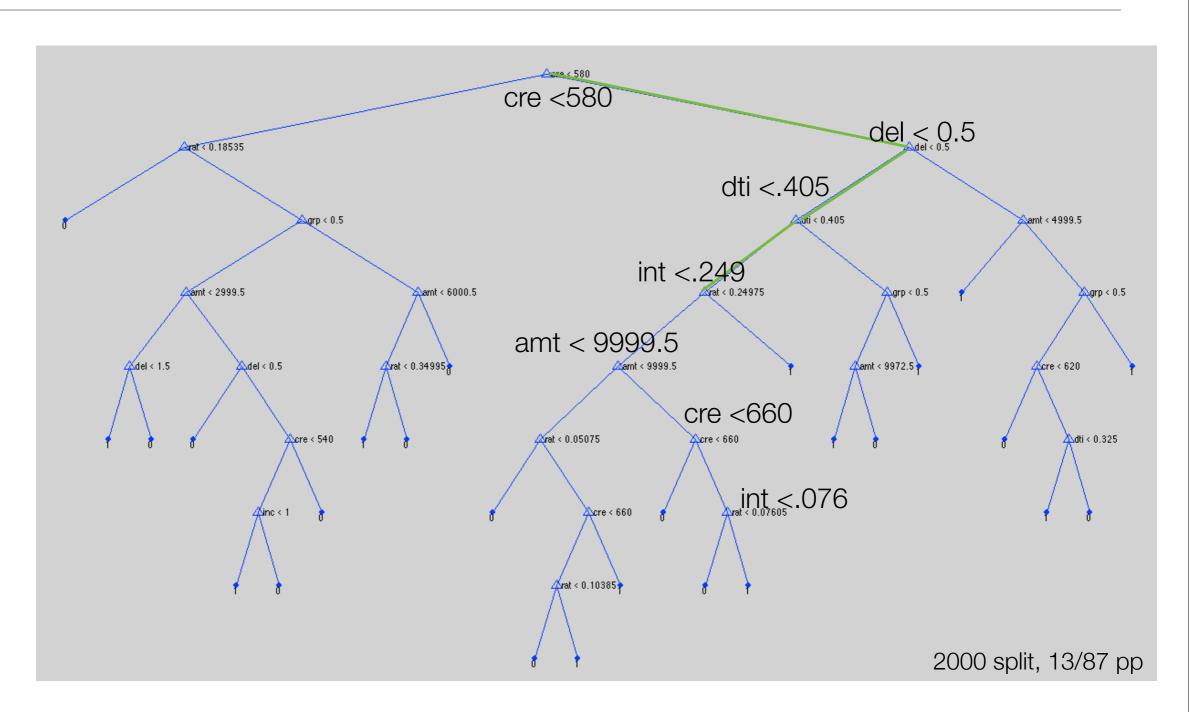
- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan



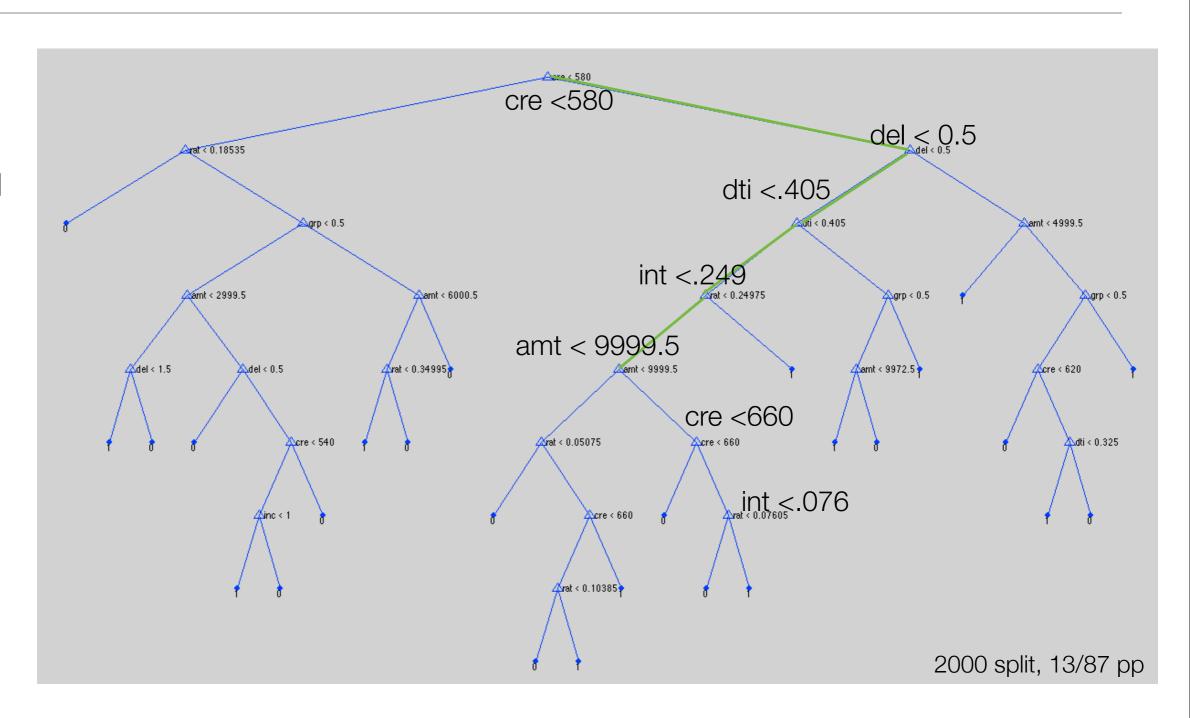
- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan



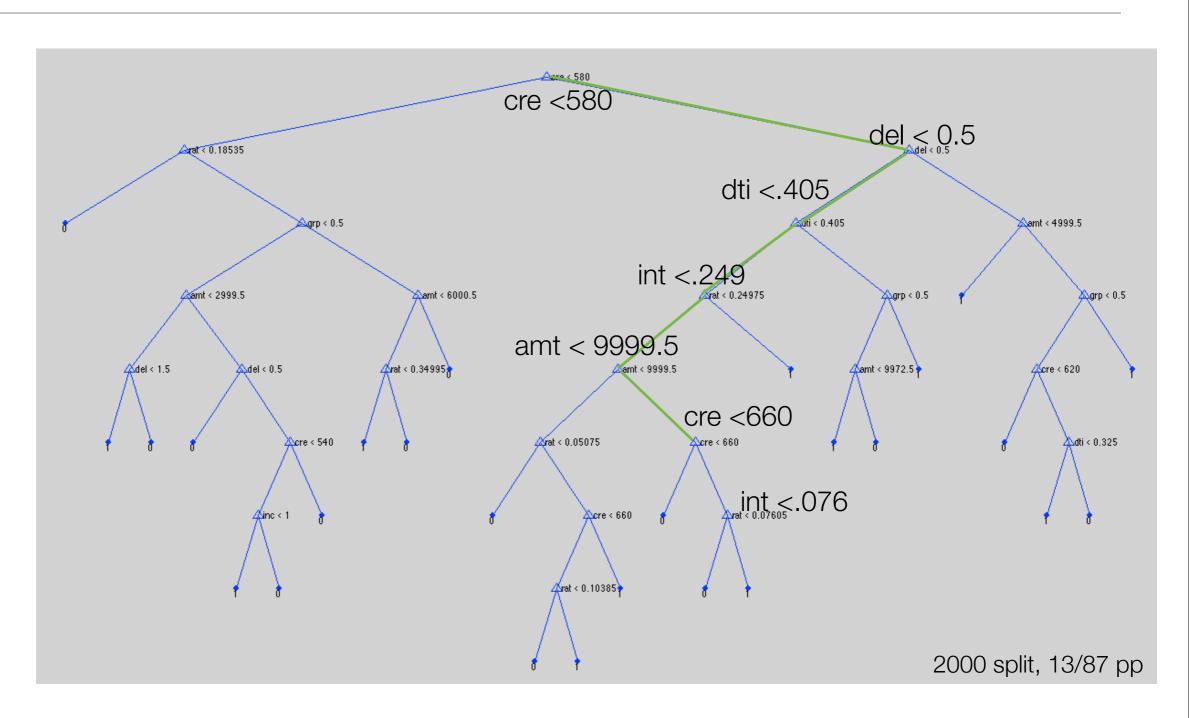
- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan



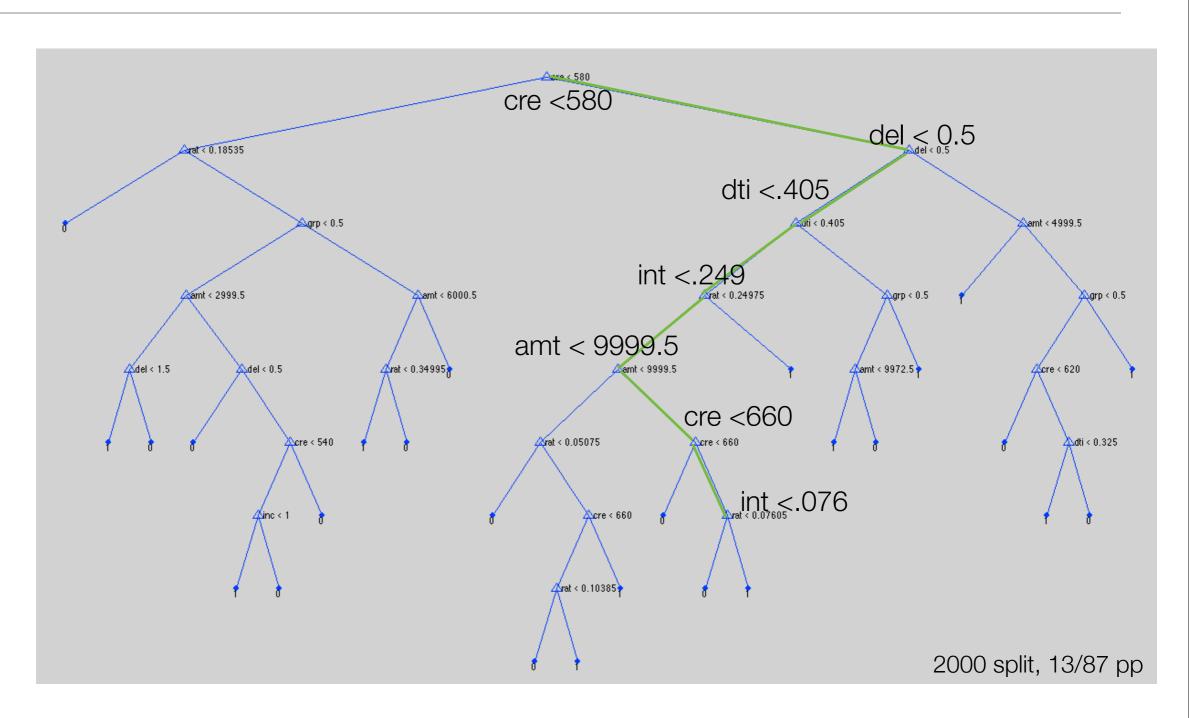
- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan



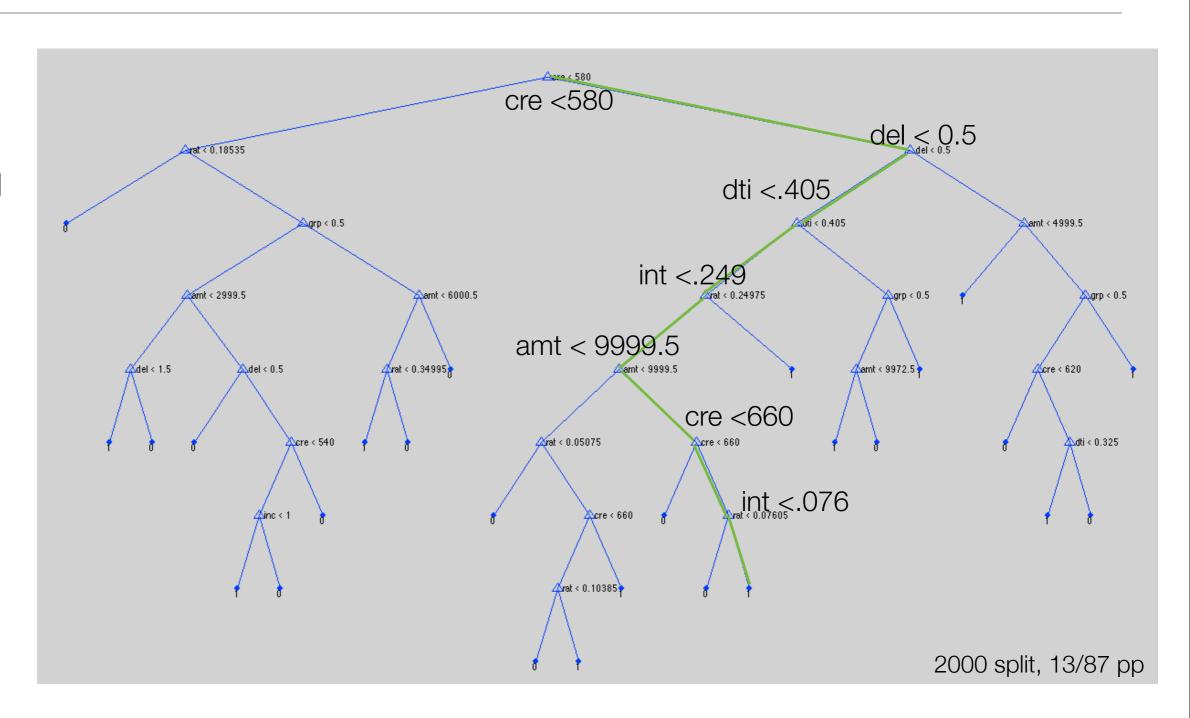
- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan



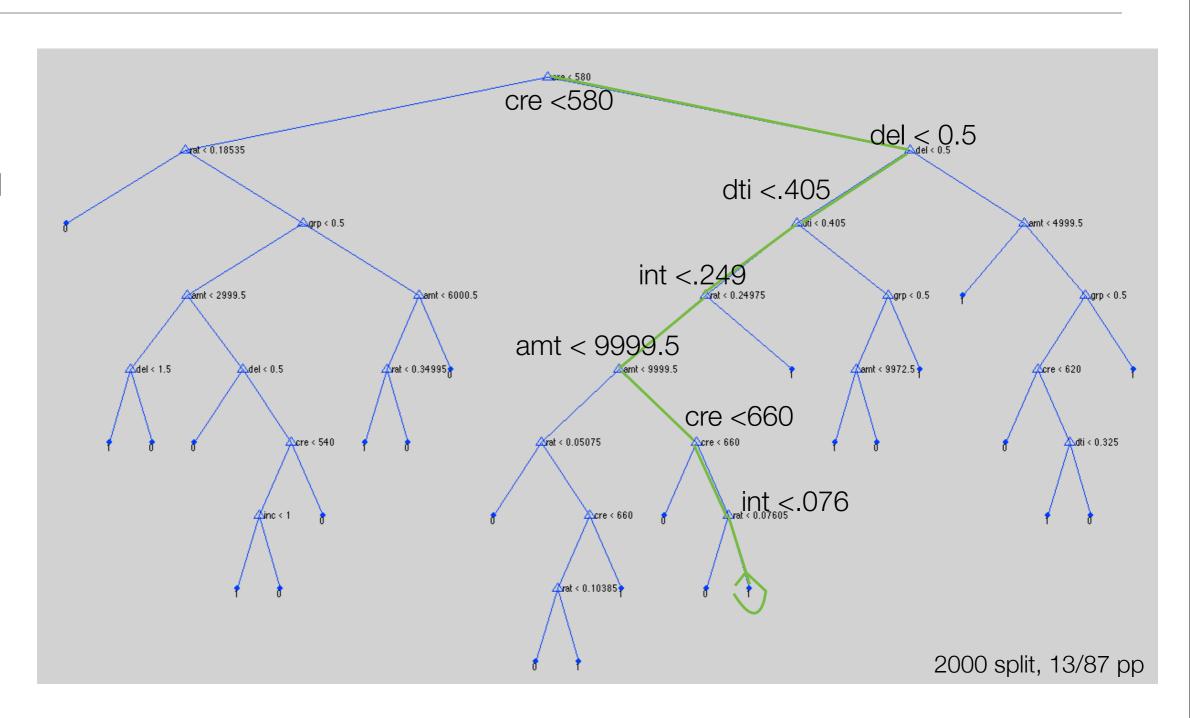
- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan



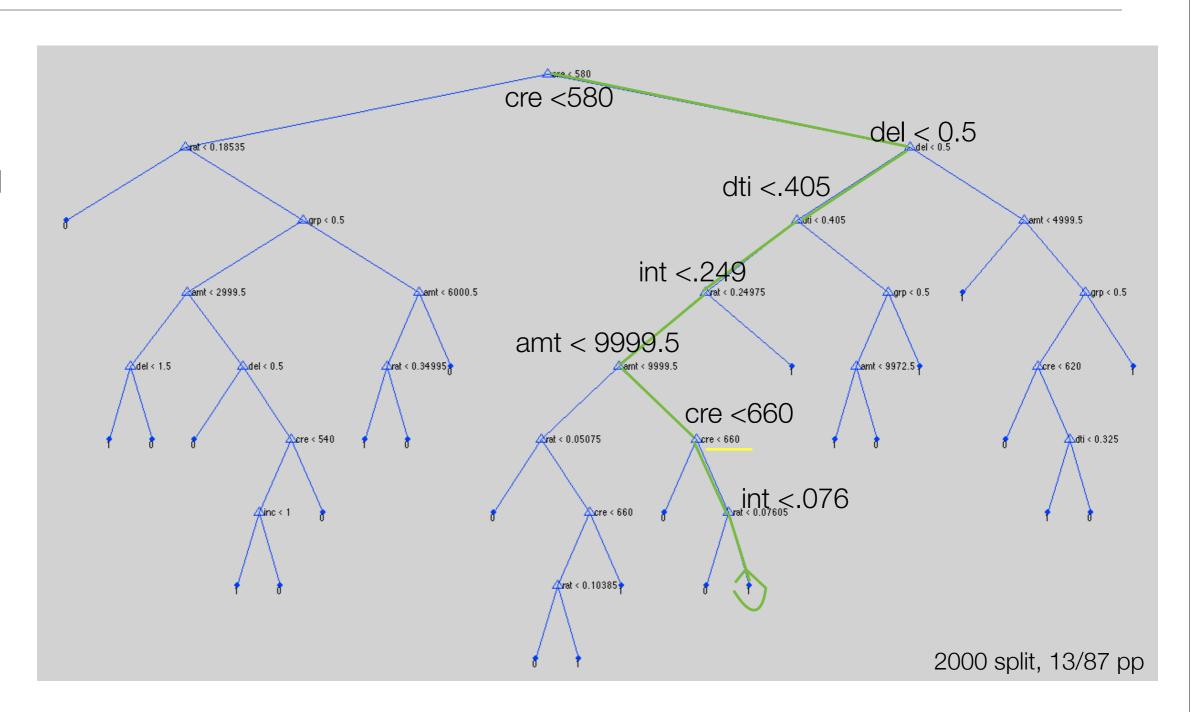
- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan



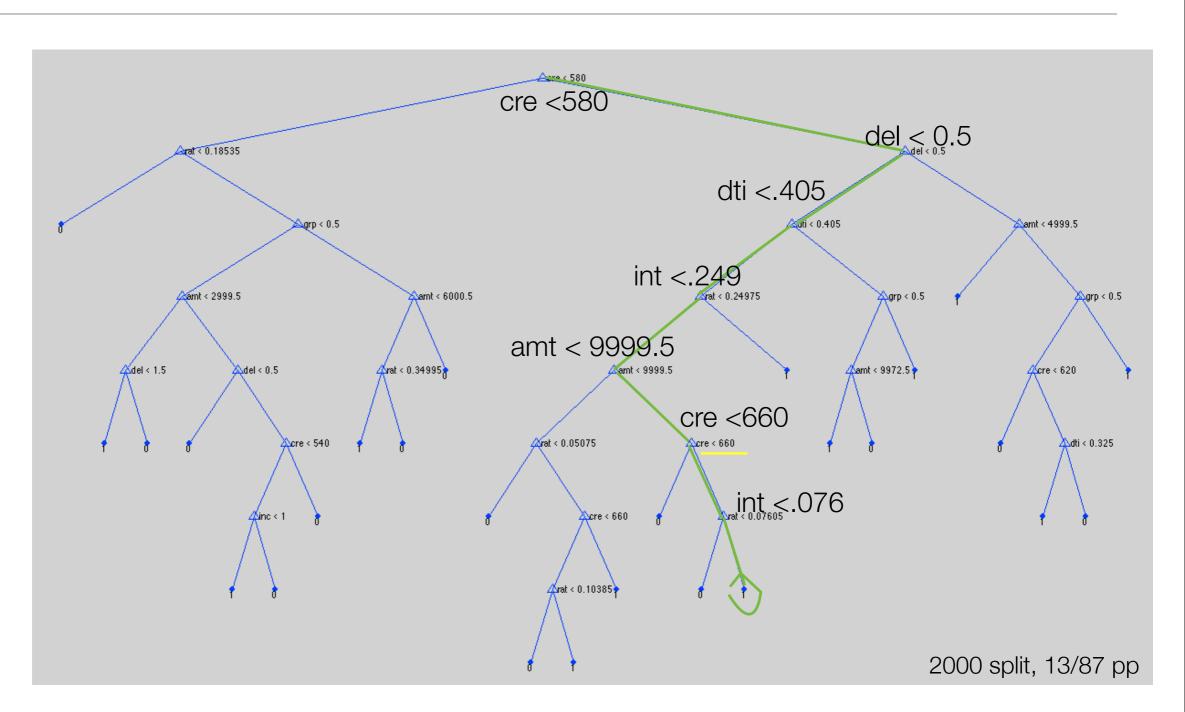
- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan



- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan

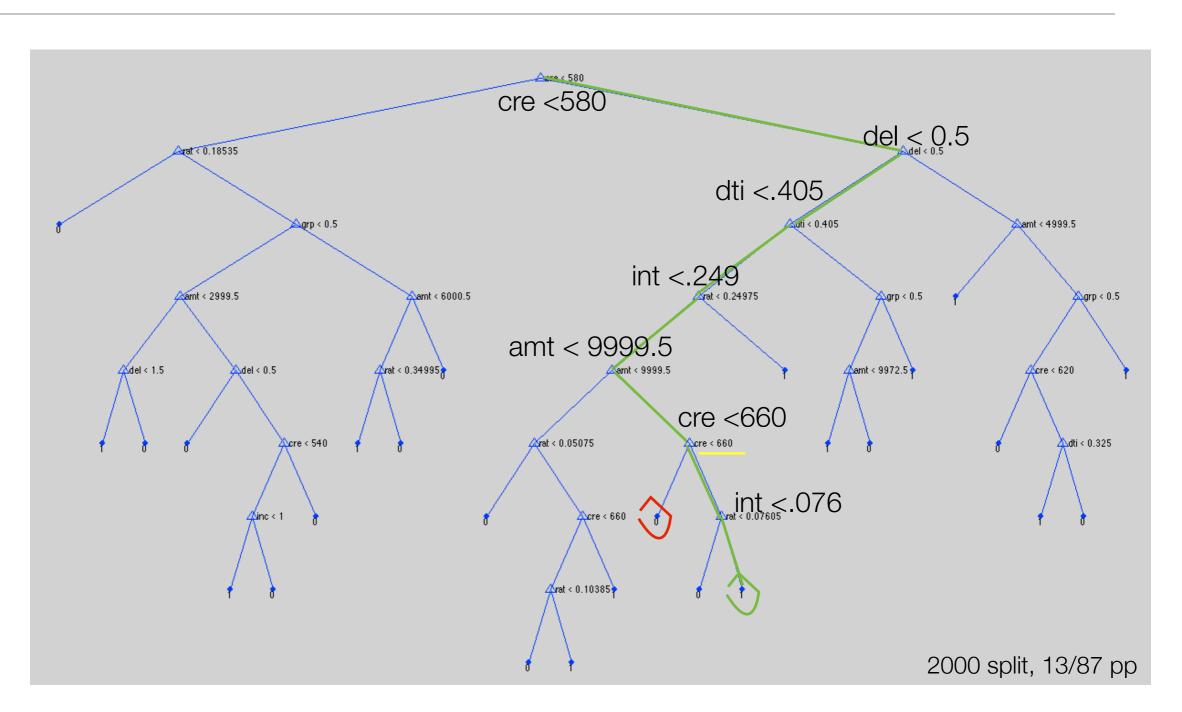


- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan



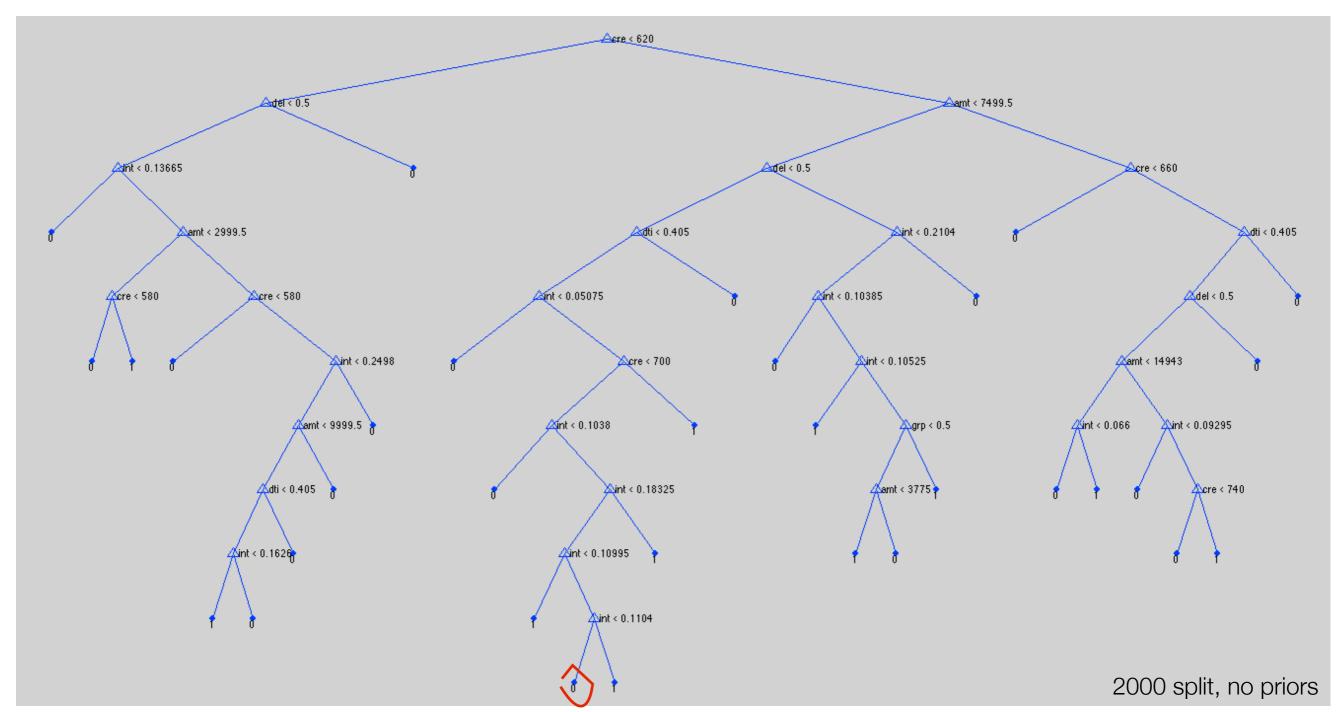
but not with C credit!

- B credit
- no delinq
- DTI 10%
- 11% interest
- \$1500 loan



but not with C credit!

priors matter! (same profile, without priors, predicts no loan)



same profile, predicts no loan

DT could be used to help borrowers set loan amount, increase loan conversion for prosper



C credit, DTI = 10%, 1 current delinquency, needs \$6000

tree predicts

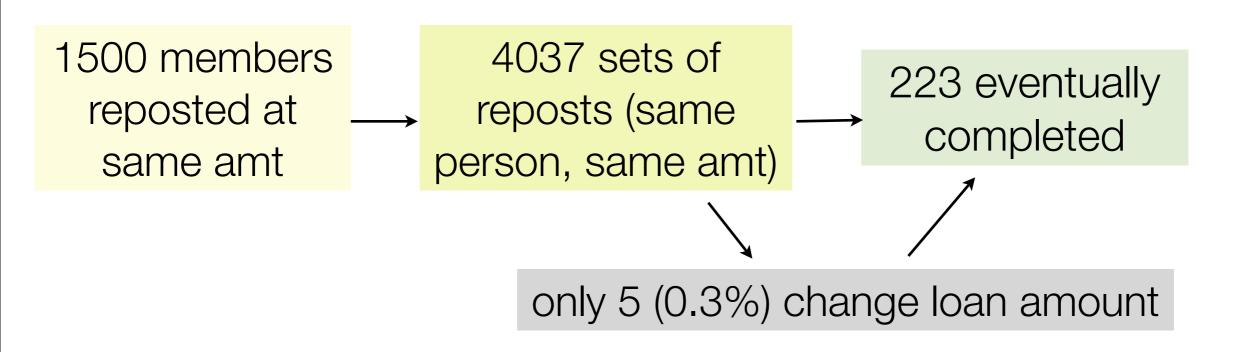
(borrowers with same profile)

56% no loan average amt listed: \$4625

44% loan average amt listed: \$3700

advice: request lower amount

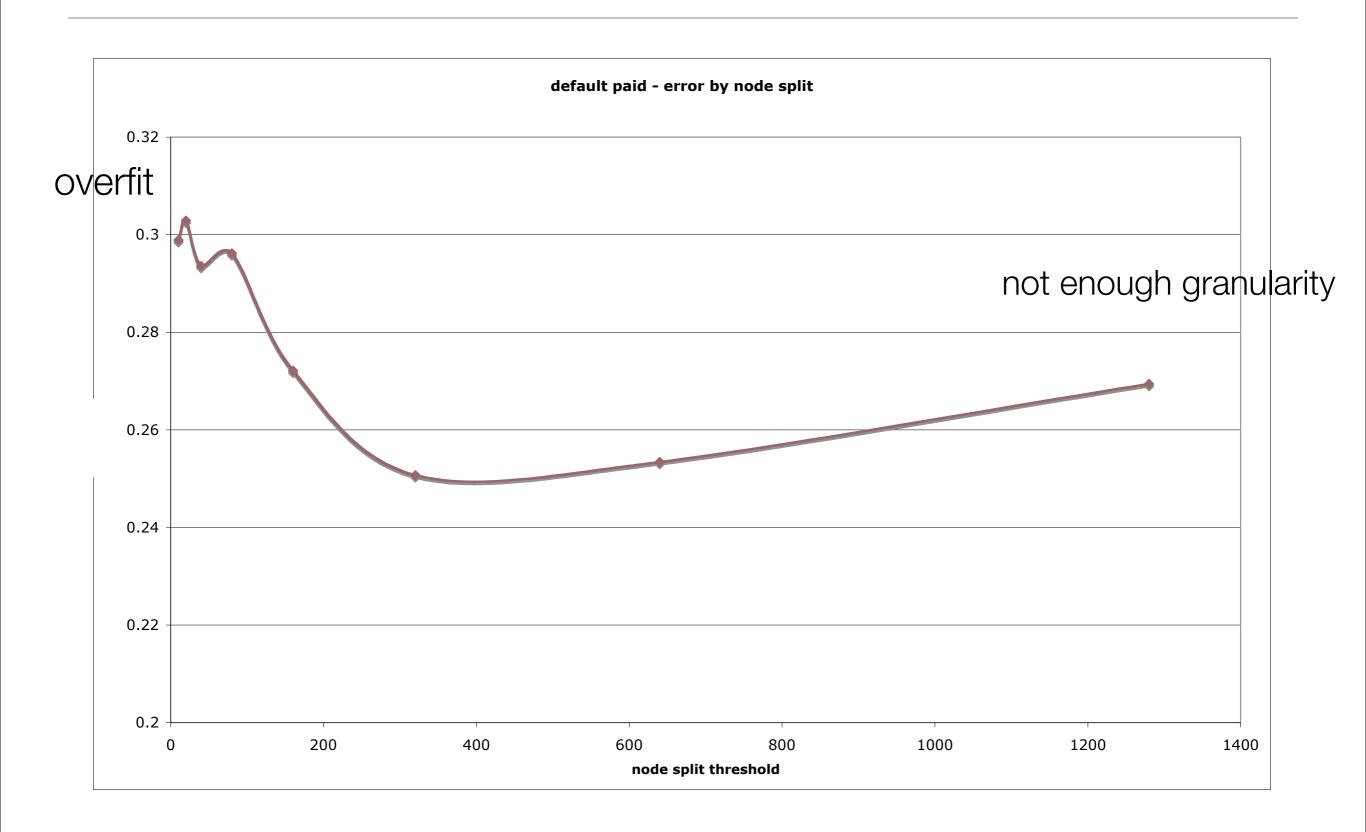
(analysis of reposted loans)



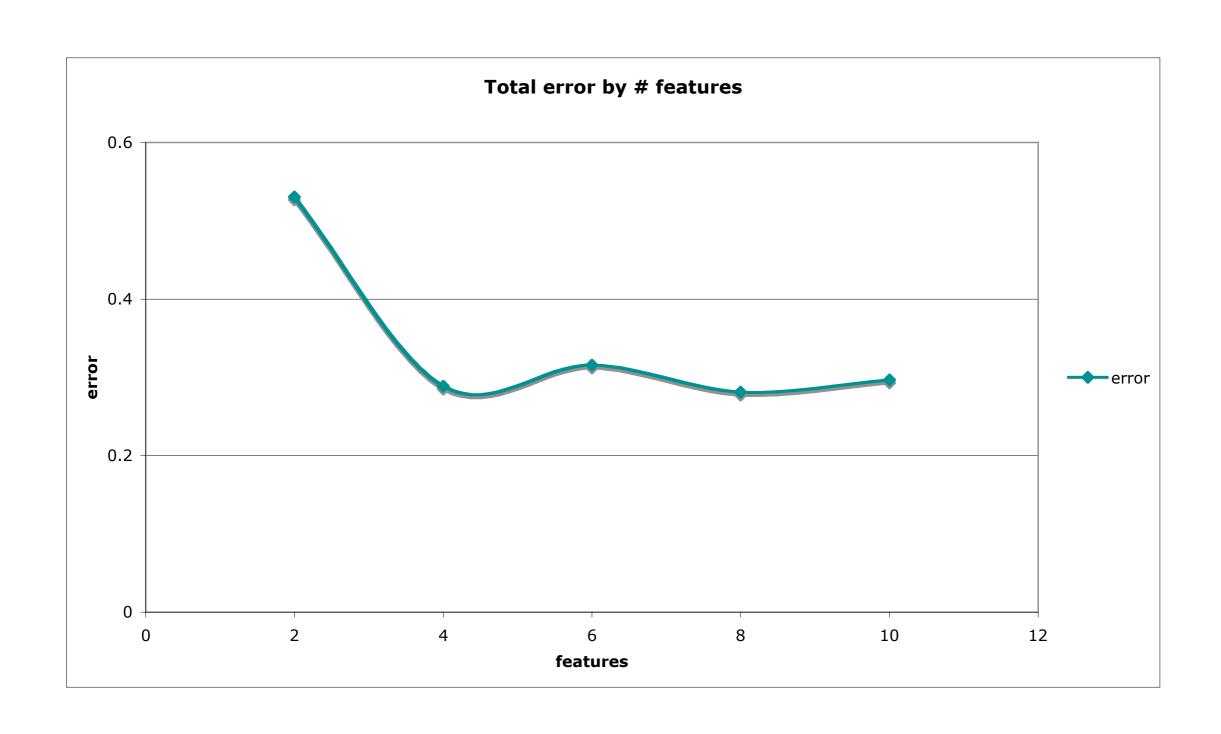
analogous process for default / paid

- Priors closer to 50/50
- Lender can use DT to identify conditional probability of default given Profile X
- Important for Prosper: keep tabs on loans with high default risk

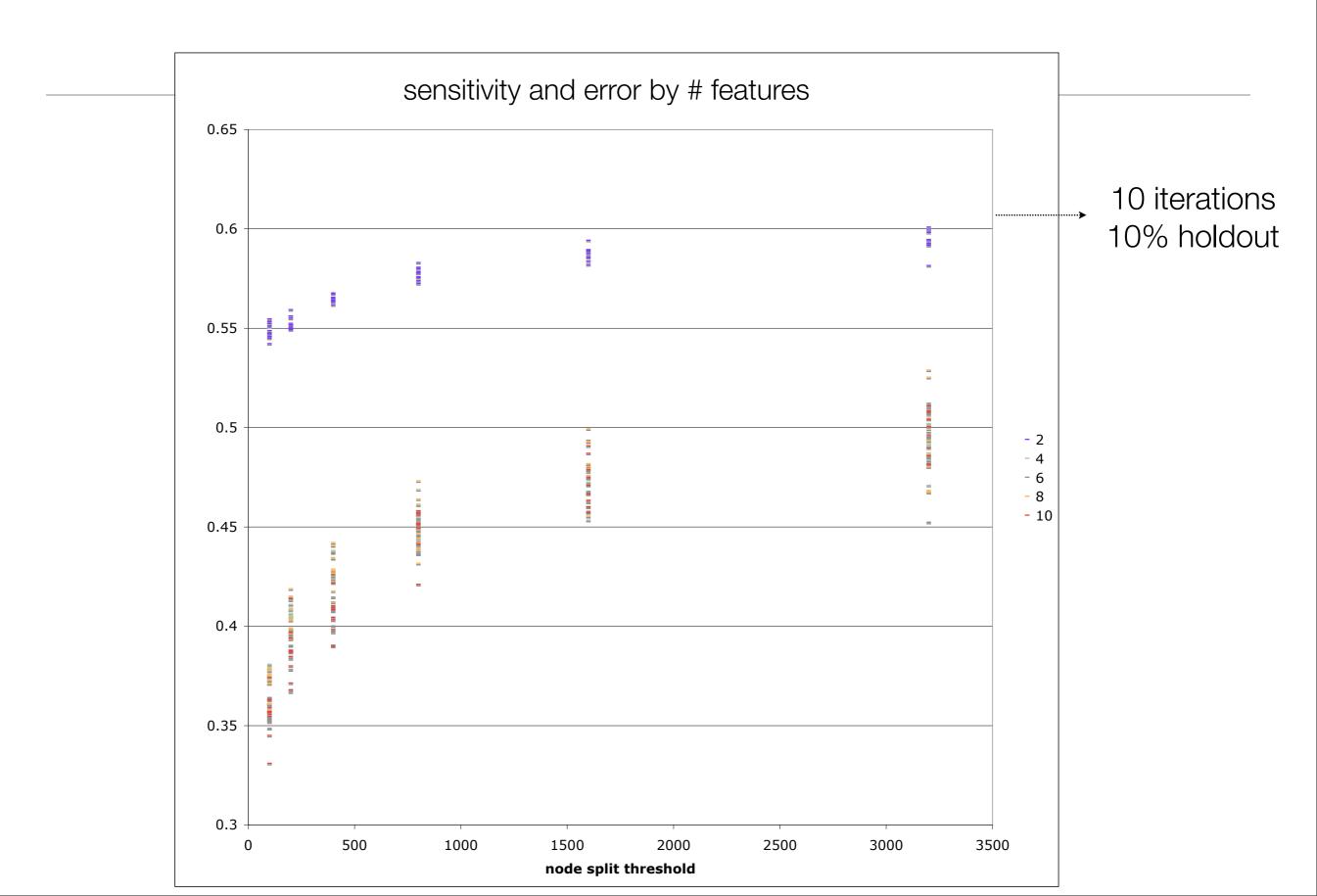
optimal pruning level (default paid tree, 6 features)



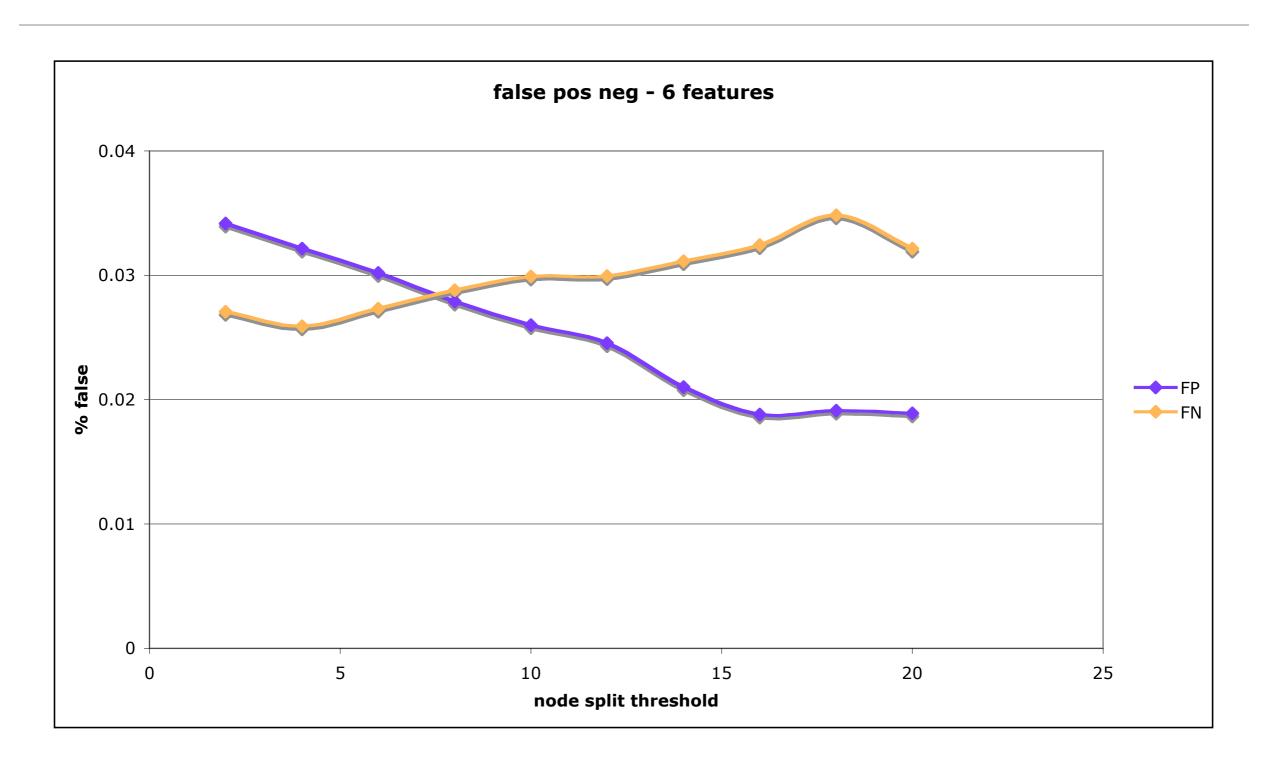
as features are added, error rate down (loan / no loan, 200 split threshold)



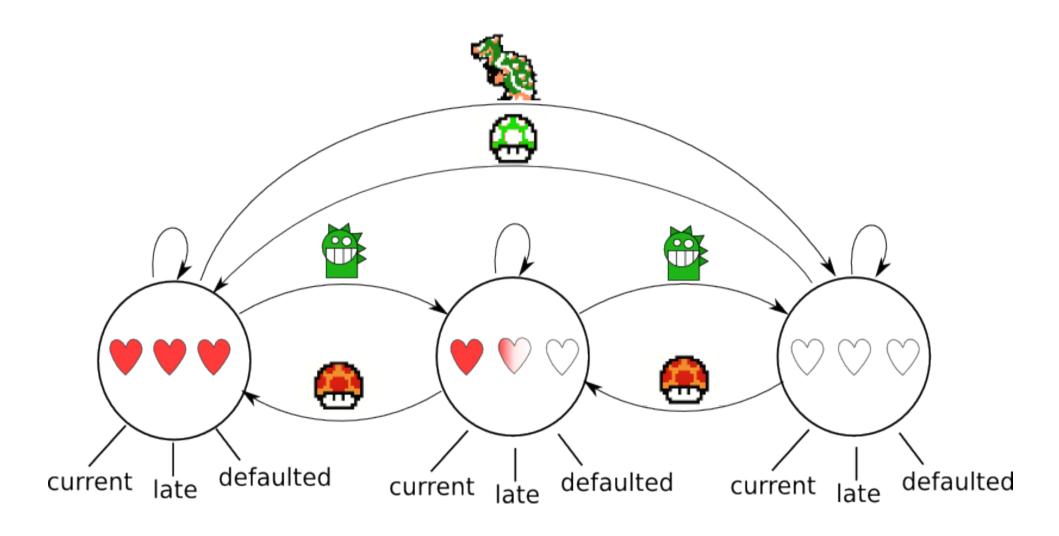
adding features, reducing split level decreases error, sensitivity



FP down, FN up as node split threshold increases (loan no loan, 6 features)

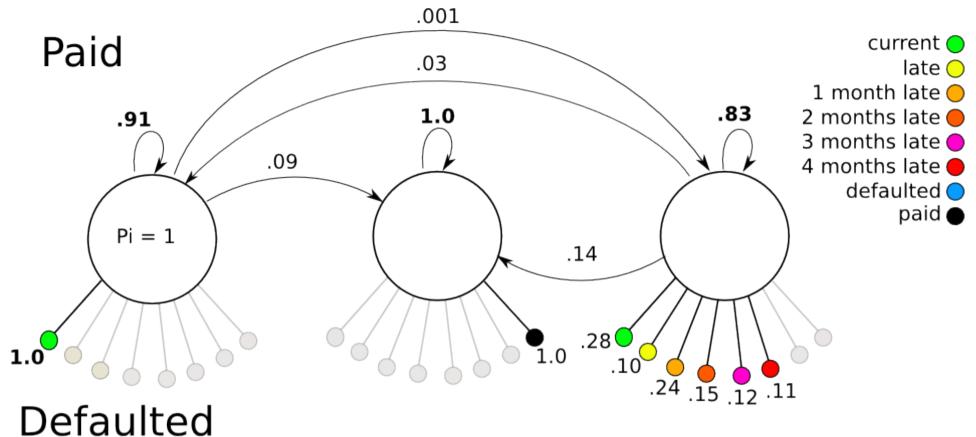


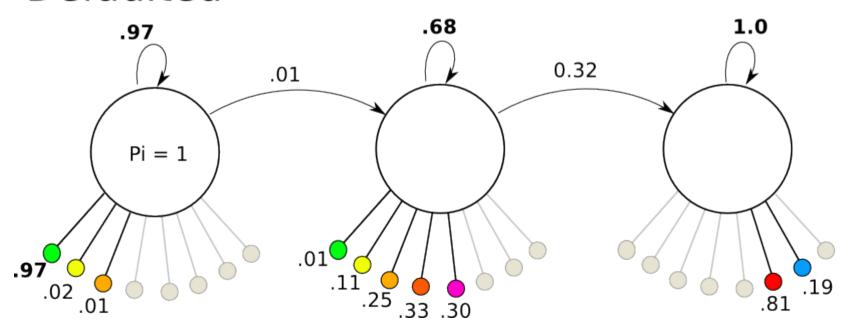
Loan Performance HMMs

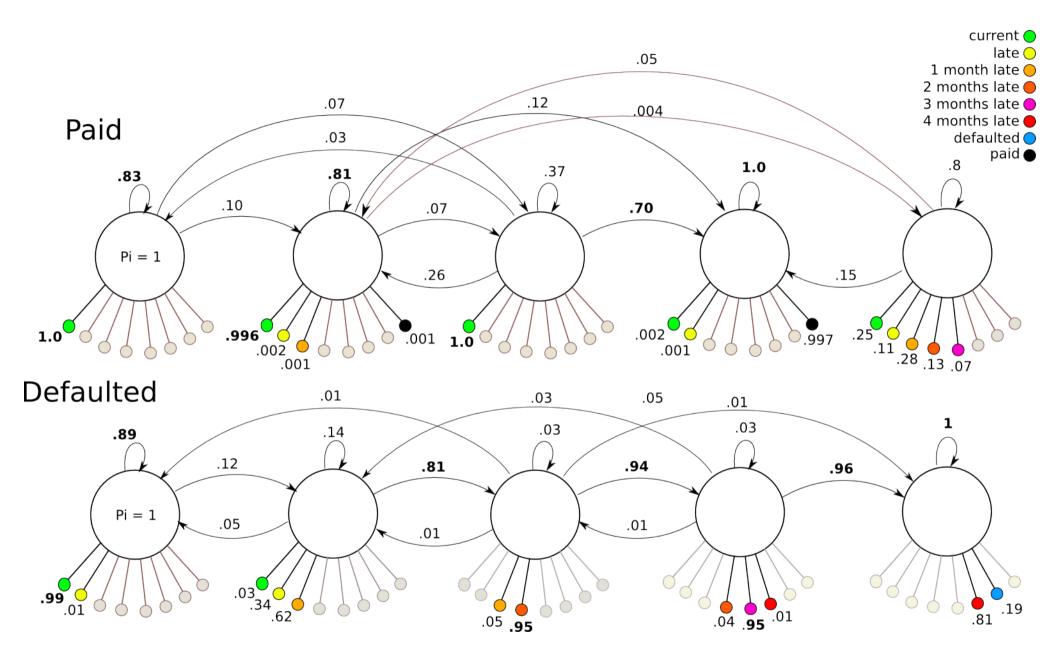


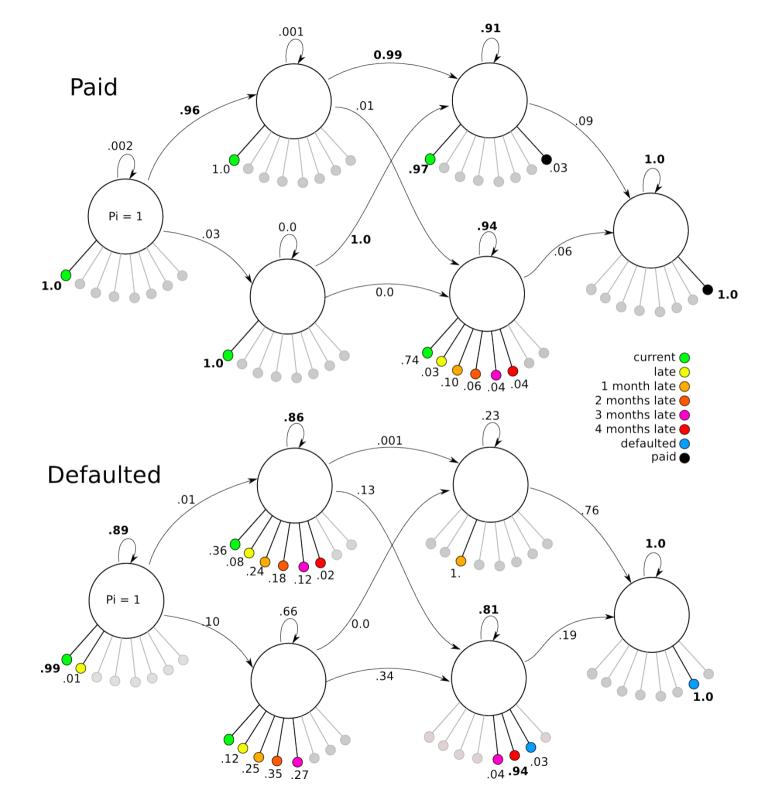
HMM Performance

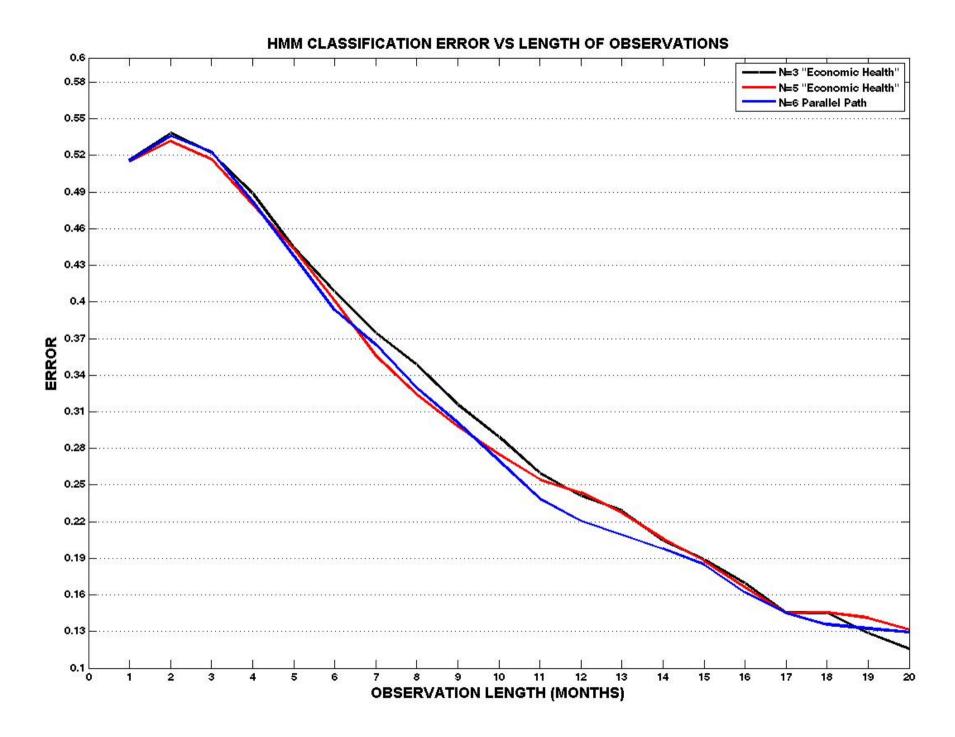
- Training 70%, testing 30% of loan performance data
 - Paid: 4365 sequences
 - Default: 3793 sequences
 - All sequences of varying length
- Model was verified and then tested











HMM Results

- Why build a model for someone who will default?
 - Short term \ long term visibility of loan performance is important
 - Default behavior potentially mimics fraud
 - A lender may be well aware of loan performance but what about Prosper?
 - Improved customer service easier to monitor high risk loans, early contact of collection agency
 - Default Performance may mimic fraud (Prosper has problems with this).

Improving HMM Performance

- Retrain specifying pi, currently start at state 1
- Create a single model of financial health
 - Train model using both paid / default data
 - Use Viterbi algorithm to estimate "proximity" to hidden state that best characterizes defaults (easy)
- Hierarchical HMM (complex)
 - Advantage is that HMMs can emit sequences of observations
 - A way to reduce error in early stages?
 - Reading
 - S. Fine, Y. Singer and N. Tishby, "The Hierarchical Hidden Markov Model: Analysis and Applications", Machine Learning, vol. 32, p. 41-62, 1998

not all groups are created equal: refining social features

 Some groups have much higher funding rates than others (queried tier-2) description, sort by category, % members funded by individual group)

highest % funded

Albuquerque

Aviation

Greece

Oil & Gas

Opthamology

Poverty Relief

Rhode Island

Rugby

Seattle

Space

Theatre

West Virginia

Veterinary

best funded, popular groups

Florida

Extended Families

Research & Analysis

Massachusetts

Travel

Accounting

Pennsylvania

Software

Financial Planning

Mortgage

Small & Medium Business

Investment Management

Family Owned

Virginia

Financial Consultants

Education & Training

Large Families

Catholic

lowest % funded

Mutual Funds

Neo- & Reform Hindus

Adoption Agencies

Air Quality

Amateur Beading

Big East Conference

Chemical

Construction

Deist DJs

Equipment & Tools

Estimating Fiction

Gambling

Glass Gliding Jewelry

Kentucky Kung Fu

Law Firms

Oceania Poetry

Printmaking

Recycling Refugees

Republicans

Security

Senior Citizens **Sporting Goods**

Structural

Surety Bonds Symphony

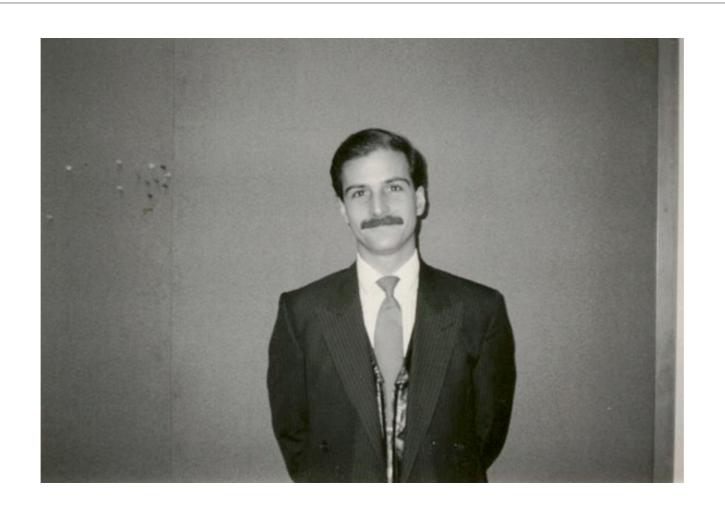
Thailand

Trading Cards United Kingdom

Utilities Yoga

Youth

Using Amazon Mechanical Turk to classify images



a human classification experiment

- Unlike banks, prosper lenders can weigh more that "just the numbers"
- Banks seek ROI; prosper lenders may have other motives (e.g. social good)
- Prosper lenders lack complex risk algorithms of banks
- Many borrowers may meet a lender's baseline criteria (e.g. FICO > 600) ... social criteria and profile assessment needed to decide how to allocate funds
- Holistic assessment of borrower profile: necessary and natural

does a borrower seem "trustworthy"?

- Goal 1: image classification
- Goal 2: assessment of "trustworthiness"
- Does "trustworthiness" correlate with getting a loan?
- Here, only pilot of methodology and analysis
- Follow-up could use humans to train classifier or create feature vector

amazon mechanical turk

Tag this image

Guidelines:

- Check the best description of what's in the image
 Check the best answer to the question: "Does this person (or the person who posted the image) look trustworthy?"
- · Your answer to the first question will be calibrated against others to ensure correct tagging

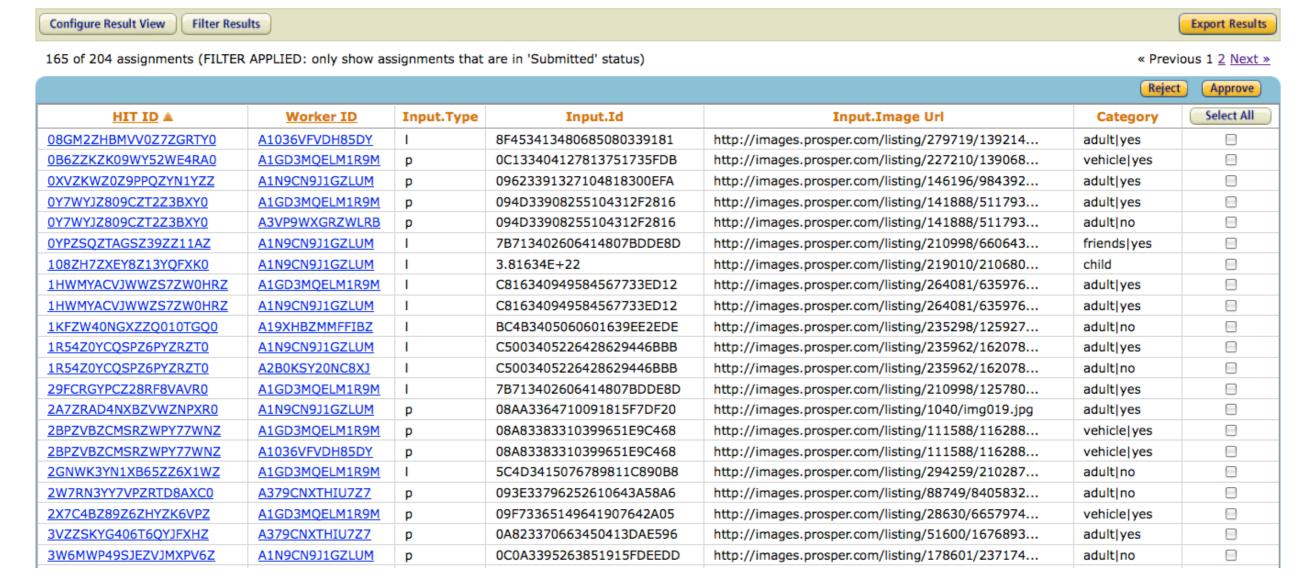
Image:



What's in this picture?			
○ Adult	Child/children	Friends/Family	○ Landscape
Animal	O House	Vehicle	Other
Does this person (or the person who posted the image) look trustworthy? Trustworthy Untrustworthy			

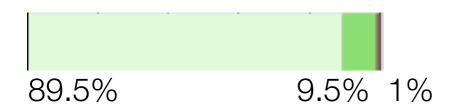
data collection

- 200 images x 2 questions x 3 workers / image (used to check for consistency)
- 50% images from unfulfilled listings; 50% from paid-off loans



consistency was good, especially for categorization

• CATEGORIES: confusion from label choice; 9.5% between children/family



• (1% disagreed on how to categorize e.g. a vehicle + people)





Multiple opinions good as fuzzy categorization?

trust rating requires clarification

• TRUST: 11% disagreement, both contextual and subjective

Lack of context





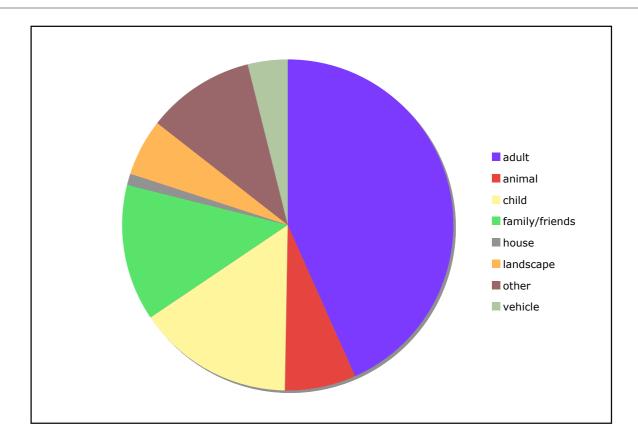
• Blurry photo, real distrust?



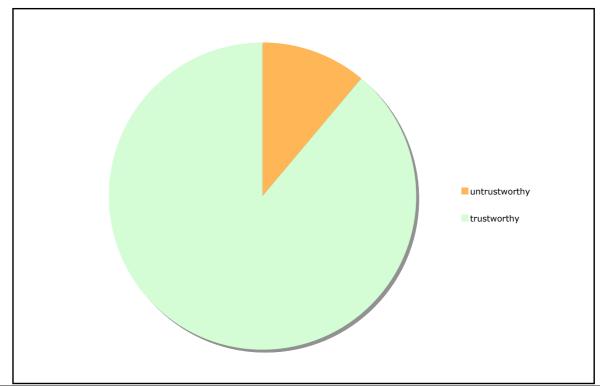


categories and trust: 71% of photos have people

categories



trust



no correlation between getting a loan and trust tag

- an image tagged "untrustworthy" was just as likely to have received & paid a loan as to have listed with no loan (no statistical difference)
- would adding contextualization (listing description) or refining the question phrasing help classification?
- Research question: how independent is judgment of "trustworthiness" from the stories built from contextual information (credit score, loan purpose), especially for quick (~8 seconds / photo) decisions?

human classification: analysis

- human-augmented classification can work: consistency was high
- experiment design is important: vague questions yield vague results
- future work could collect larger sample; use as a feature vector
- also, text / spelling analysis