



	h

<i>Output = Predicted</i>	Reply Time (neares	t hour)		Input = <email_fea< th=""><th>atures></th><th></th><th>Talaa</th></email_fea<>	atures>		Talaa
Training Size = 4,04 Test Size = 1,011	44 1) Aggregated Dataset	t	 Iakeaways: Heavy skew in our data: all models except Linear 	<i>Output = Predicted Training Size = 955 Test Size = 239</i>	Reply Time (neares	t hour)	• Po wit
	Training Accuracy	Test Accuracy	Rearession predict		(Z) Balanced Dataset	Tost Accuracy	• Po
Linear Regression	0.061	0.059	the most common	Linoar Pogrossion	0 055	0 055	dat
Logistic Regression	0.368	0.362	label	Logistic Regression	0.033	0.033	
Naïve Bayes	0.509	0.380		Naïve Baves	0.137	0.007	
Neural Network	0.322	0.319	 Dataset not linearly 	Neural Network	0.100	0.066	
Experiment #3	: Used a balanc	ed dataset wh	ere we "binarize" our	Experiment #4	Chose two mos	st promising r	model
Experiment #3 labels into class Input = <email_fea Output = 1 if reply t</email_fea 	: Used a balanc ses of reply time atures> ime > 30 mins, 0 oth	ed dataset wh s < 1/2 hr and	ere we "binarize" our reply times >= 1/2 hr	Experiment #4 with same data Input = <email_fea Output = 1 if reply t</email_fea 	<i>Chose two mosset as prior, but atures</i> <i>time > 30 mins, 0 oth</i>	st promising r with improved	model d featu Take
Experiment #3 labels into class Input = <email_fea Output = 1 if reply to Training Size = 2,88</email_fea 	<i>: Used a balanc</i> ses of reply time atures> ime > 30 mins, 0 oth 96	ed dataset wh s < 1/2 hr and	ere we "binarize" our reply times >= 1/2 hr Takeaways:	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,8</email_fea 	<i>Chose two mosset as prior, but atures</i> <i>ime > 30 mins, 0 oth</i>	st promising r with improved	model d featu Takea
Experiment #3 labels into class Input = <email_fea Output = 1 if reply t Training Size = 2,89 Test Size = 724</email_fea 	: Used a balanc Ses of reply time atures> ime > 30 mins, 0 oth 96	ed dataset wh s < 1/2 hr and	ere we "binarize" our reply times >= 1/2 hr Takeaways: • Slightly better than random performance	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,8 Test Size = 724</email_fea 	<i>Chose two mosset as prior, but atures</i> <i>ime</i> > 30 mins, 0 oth	st promising r with improved	model d featu Take a • He
Experiment #3 labels into class Input = <email_fea Output = 1 if reply t Training Size = 2,89 Test Size = 724</email_fea 	: Used a balanc Ses of reply time atures> ime > 30 mins, 0 oth 96 (3) Binarized Dataset Training Accuracy	ed dataset wh s < 1/2 hr and herwise Test Accuracy	ere we "binarize" our reply times >= 1/2 hr Takeaways: • Slightly better than random performance on all models.	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,8 Test Size = 724</email_fea 	<i>Chose two mosset as prior, but atures</i> <i>ime</i> > 30 mins, 0 oth	st promising r with improved	model d featu • He neu ove
Experiment #3 labels into class Input = <email_fea Output = 1 if reply t Training Size = 2,89 Test Size = 724</email_fea 	: Used a balanc ses of reply time atures> ime > 30 mins, 0 oth 96 (3) Binarized Dataset Training Accuracy 0.536	ed dataset wh s < 1/2 hr and herwise Test Accuracy 0.526	ere we "binarize" our reply times >= 1/2 hr Takeaways: • Slightly better than random performance on all models.	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,83 Test Size = 724 (4) Binariz</email_fea 	Chose two mos set as prior, but atures> ime > 30 mins, 0 oth 96	st promising i with improved	model d featu • He neu ove
Experiment #3 labels into class Input = <email_fea Output = 1 if reply t Training Size = 2,89 Test Size = 724</email_fea 	: Used a balanc ses of reply time atures> ime > 30 mins, 0 oth 96 (3) Binarized Dataset Training Accuracy 0.536 0.543	ed dataset wh s < 1/2 hr and nerwise <u>Test Accuracy</u> 0.526 0.531	ere we "binarize" our reply times >= 1/2 hr Takeaways: • Slightly better than random performance on all models.	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,83 Test Size = 724 (4) Binariz</email_fea 	Chose two mos set as prior, but atures> ime > 30 mins, 0 oth 96 ed Dataset + Improved Training Accuracy	st promising i with improved Features Test Accuracy	model d featu Takea • He neu ove
Experiment #3 labels into class lnput = <email_fea Output = 1 if reply t Training Size = 2,88 Test Size = 724</email_fea 	: Used a balanc ses of reply time atures> ime > 30 mins, 0 oth 96 (3) Binarized Dataset (3) Binarized Dataset 0.536 0.543 0.543	<i>ed dataset wh</i> <i>s < 1/2 hr and</i> <i>nerwise</i> 0.526 0.531 0.500	ere we "binarize" our reply times >= 1/2 hr Takeaways: • Slightly better than random performance on all models. • SVM and NN perform best. Both are	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,8 Test Size = 724 (4) Binariz</email_fea 	Chose two mos set as prior, but atures> ime > 30 mins, 0 oth 96 ced Dataset + Improved Training Accuracy 0.773	st promising i with improved ferwise Features Test Accuracy 0.560	model d featu • He neu ove
Experiment #3 labels into classlaput = <email_fea< td="">Output = 1 if reply tTraining Size = 2,88Test Size = 724Linear RegressionLogistic RegressionNaïve BayesSVM</email_fea<>	: Used a balanc ses of reply time atures> ime > 30 mins, 0 oth 96 (3) Binarized Dataset (3) Binarized Dataset (3) Binarized Dataset 0.536 0.536 0.543 0.502 0.871	<i>ed dataset wh</i> <i>s < 1/2 hr and</i> <i>nerwise</i> 0.526 0.531 0.500 0.558	ere we "binarize" our reply times >= 1/2 hr Takeaways: • Slightly better than random performance on all models. • SVM and NN perform best. Both are overfitting	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,8 Test Size = 724 (4) Binariz SVM Neural Network</email_fea 	Chose two mos set as prior, but atures> ime > 30 mins, 0 oth 96 ced Dataset + Improved Training Accuracy 0.773 0.828	st promising i with improved erwise Features Test Accuracy 0.560 0.570	model d featu Takea • He neu ove

Training Size = $4,044$ Test Size = $1,011$		 Takeaways: Heavy skew in our 	Input = <email_features> Output = Predicted Reply Time (nearest hour) Training Size = 955</email_features>			Take a • Po	
			data: all models	Test Size = 239			wit
(1	1) Aggregated Datase	t	except Linear		(2) Balanced Dataset		
	Training Accuracy	Test Accuracy	Regression predict		Training Accuracy	Test Accuracy	• Po
Linear Regression	0.061	0.059	the most common	Linear Regression	0.055	0.055	da
Logistic Regression	0.368	0.362	label	Logistic Regression	0.137	0.097	
Naïve Bayes	0.509	0.380		Naïve Bayes	0.509	0.008	
Neural Network	0.322	0.319	 Dataset not linearly 	Neural Network	0.100	0.066	
Experiment #3	: Used a balanc	ed dataset wh	ere we "binarize" our	Experiment #4	Chose two mo	st promising I	model
Experiment #3 labels into class Input = <email_fea Output = 1 if reply the</email_fea 	<i>: Used a balanc</i> ses of reply time tures> ime > 30 mins, 0 oth	ed dataset wh s < 1/2 hr and	ere we "binarize" our reply times >= 1/2 hr Takeawavs:	Experiment #4 with same data Input = <email_fea Output = 1 if reply t</email_fea 	<i>Chose two mosset as prior, but atures So mins, 0 oth</i>	st promising i with improved	model d feate
Experiment #3 labels into class Input = <email_fea Output = 1 if reply the Training Size = 2,89</email_fea 	<i>: Used a balanc</i> ses of reply time tures> ime > 30 mins, 0 oth 96	ed dataset wh s < 1/2 hr and	ere we "binarize" our reply times >= 1/2 hr Takeaways:	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,83</email_fea 	<i>Chose two mosset as prior, but atures</i> <i>ime > 30 mins, 0 oth</i>	st promising i with improved	model d feate
Experiment #3 labels into class Input = <email_fea Output = 1 if reply to Training Size = 2,89 Test Size = 724</email_fea 	: Used a balance ses of reply time tures> ime > 30 mins, 0 oth 26	ed dataset wh s < 1/2 hr and herwise	ere we "binarize" our reply times >= 1/2 hr Takeaways: • Slightly better than random performance	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,83 Test Size = 724</email_fea 	<i>Chose two mosset as prior, but atures</i> <i>ime</i> > 30 <i>mins,</i> 0 <i>oth</i>	st promising i with improved	model d feate • He
Experiment #3 labels into class Input = <email_fea Output = 1 if reply to Training Size = 2,89 Test Size = 724</email_fea 	: Used a balanc ses of reply time tures> ime > 30 mins, 0 oth 26 (3) Binarized Dataset Training Accuracy	ed dataset wh s < 1/2 hr and herwise	ere we "binarize" our reply times >= 1/2 hr Takeaways: • Slightly better than random performance on all models.	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,8 Test Size = 724</email_fea 	<i>Chose two mosset as prior, but atures</i> <i>ime</i> > 30 mins, 0 oth	st promising i with improved	model d feate • He net
Experiment #3 labels into class Input = <email_fea Output = 1 if reply to Training Size = 2,89 Test Size = 724</email_fea 	: Used a balance ses of reply time tures> ime > 30 mins, 0 oth 26 (3) Binarized Dataset Training Accuracy 0.536	ed dataset wh os < 1/2 hr and nerwise Test Accuracy 0.526	ere we "binarize" our reply times >= 1/2 hr Takeaways: • Slightly better than random performance on all models.	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,85 Test Size = 724 (4) Binariz</email_fea 	Chose two mo set as prior, but atures> ime > 30 mins, 0 oth 96	st promising i with improved	model d feate • He net ove
Experiment #3 labels into class Input = <email_fea Output = 1 if reply to Training Size = 2,89 Test Size = 724</email_fea 	: Used a balance ses of reply time tures> ime > 30 mins, 0 oth 26 (3) Binarized Dataset Training Accuracy 0.536 0.543	ed dataset wh os < 1/2 hr and nerwise <u>Test Accuracy</u> 0.526 0.531	ere we "binarize" our reply times >= 1/2 hr Takeaways: • Slightly better than random performance on all models. • SVM and NN perform	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,8 Test Size = 724 (4) Binariz</email_fea 	Chose two mo set as prior, but atures> ime > 30 mins, 0 oth 96 ed Dataset + Improved Training Accuracy	st promising i with improved ferwise	model d featu Takea • He neu ove
Experiment #3 Iabels into class Input = <email_fea Output = 1 if reply th Training Size = 2,89 Test Size = 724</email_fea 	: Used a balanc ses of reply time tures> ime > 30 mins, 0 oth 06 (3) Binarized Dataset Training Accuracy 0.536 0.543 0.502	ed dataset wh os < 1/2 hr and nerwise 0.526 0.531 0.500	ere we "binarize" our reply times >= 1/2 hr Takeaways: • Slightly better than random performance on all models. • SVM and NN perform best. Both are	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,8 Test Size = 724 (4) Binariz</email_fea 	Chose two mo set as prior, but atures> ime > 30 mins, 0 oth 96 ed Dataset + Improved Training Accuracy 0.773	st promising i with improved erwise Features Test Accuracy 0.560	model d featu • He net ove
Experiment #3 labels into classlaput = <email_fea< td="">Output = 1 if reply toTraining Size = 2,88Test Size = 724Linear RegressionLogistic RegressionNaïve BayesSVM</email_fea<>	: Used a balanc ses of reply time tures> ime > 30 mins, 0 oth 06 (3) Binarized Dataset Training Accuracy 0.536 0.543 0.502 0.871	ed dataset wh es < 1/2 hr and nerwise 0.526 0.531 0.500 0.558	 ere we "binarize" our reply times >= 1/2 hr Takeaways: Slightly better than random performance on all models. SVM and NN perform best. Both are overfitting 	Experiment #4 with same data Input = <email_fea Output = 1 if reply t Training Size = 2,83 Test Size = 724 (4) Binariz</email_fea 	Chose two mos set as prior, but atures> ime > 30 mins, 0 oth 96 ced Dataset + Improved Training Accuracy 0.773 0.828	st promising i with improved erwise Features Test Accuracy 0.560 0.570	model d featu • He net ove

Predicting Expected Email Response Times

Laura Cruz-Albrecht (Icruzalb@stanford.edu) and Kevin Khieu (kkhieu@stanford.edu) CS229: Machine Learning, Fall 2017

iments and Results

Features

Raw Input Features: In deciding raw input features, we chose features that we hypothesized could influence users into replying more/less quickly. These features are:

- Num. Recipients in "To" field
- Num. Words in Email
- Num. Recipients in "CC" field
- If Email is Reply
- Time of Day
- Day of Week
- Num. Words in Subject
- Num. "?" in Body/Subject
- If "?" Mark in Body/Subject
- If Keywords = ['response', 'please', 'can', 'urgent', 'important', 'need'] in Body/Subject

Derived Features: Our dataset was proven to not be linearly separable, so we used an SVM that mapped features with a Gaussian RBF Kernel

we cap the to 50

aways:

oor classification even th balanced dataset

stentially insufficient ita for 25 labels

Is and trained ure set

aways:

eavy overfitting by eural network, less verfitting by SVM

Our investigation indicates that the problem of email response time prediction is difficult: in the multilabel space, even after balancing the labels of our skewed initial dataset, prediction accuracy across models was low. The classification problem was still challenging when reduced to one of binary prediction among balanced classes. The low accuracies across models is understandable however: there are a lot of social factors that

Future Work

Given more time, the following is a list of future directions for this project:

- gathering more diverse data would boost the robustness of our models
- further feature extraction investigation could be done to find more indicative features
- incorporating each user's history of reply **behavior**, which would require gaining user permission, would greatly aid in our work



Voces

 $h(x) = \sum \theta_i x_i = \theta^T x$

 $h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$

Naive Bayes Binary + Multinomial **TF-IDF** Standardization

Linear Regression

Ordinary Least Squares

Logistic Regression

Maximizing Log Likelihood

Minimization via SGD

Binary + Multiclass

Support Vector Machine Gaussian RBF Kernel + Polynomial Kernel

Neural Network Multilayer Perceptron 1-2 Hidden Layers (60x30)

$\hat{y} = rgmax_{k \in \{1,\ldots,K\}} p(C_k) \prod_{i=1}^n p(x_i \mid C_k)$

Discussion

cannot be detected in the metadata of an email that impact how quickly a person chooses

to, or is able to, respond. Nevertheless, we did find that the Neural Net and SVM Models did perform at above-random levels, indicating that there is some degree of correlation between features of an email, and the likelihood of receiving a quick response within the work setting.

Sent

Email --

