*15.780 Final Project Report*
*Ankita Devasia, Deekshita Kacham, Lydia Yu*

# Table of Contents

## 1.1 Background

Our project aims to supplement work done by the Warehouse Work, Health, and Well-Being study conducted by the Sloan School of Management under Professor Erin Kelly. The purpose of this study is to to understand what it's like to work at a fulfillment center for an online retailer and to understand how changes in the workplace affect workers' experiences. It also hopes to investigate how work conditions may affect health and well-being, personal and family life, and decisions about whether to stay in the job or leave.

Honing in on one subset of the problem, our team hopes to supplement the work already done by the lab and find relationships between particular variables and the productivity of the warehouse. Upon doing some initial research, we have decided on the problem statement discussed in the next section. More information on the study itself can be found at https://mitsloan.mit.edu/programs/phd/mit-warehouse-work-study.

## 1.2 Significance

The Warehouse Worker Well-being study seeks to document ways to improve workers' experience in warehouses, a setting rapidly being affected by e-commerce and more recently by the pandemic. This project aims to investigate how work redesign – deliberate changes in work conditions – may improve the health and well-being of workers.

The pandemic has significantly increased consumers' reliance on online shopping, putting strain on many warehouse workers. In light of this, Professor Kelly's lab partnered with the e-commerce division of a national retail firm to introduce a participatory workplace intervention and its impact on employees' health and wellbeing and on key organizational outcomes. The study has just finished initial rounds of survey data collection in which workers of the various warehouses of this retailer were surveyed on topics related to their wellbeing and work conditions. The results of these surveys are currently being communicated to the managers of each of the fulfillment centers at which the surveys were sent out.

A significant measure of worker wellbeing is how likely they are to leave their jobs or be laid off. If turnover has a critical impact on the productivity of a warehouse, then this relationship could be an important target for our study because productivity is directly tied with profits and other measures of business outcomes.

## 2. Objective

Our goal is to support the Warehouse Work, Health, and Well-Being study by studying the effect of turnover on productivity within warehouses of a national online retailer and finding an optimal level of turnover that maximizes productivity. We have labor productivity data from an E-commerce retailer as well as monthly turnover data for the retailer's warehouses, which corresponds to how many workers leave each fulfillment center (either voluntarily or through lay-offs) each month.

## 3. Dataset

Through this study, we have acquired data on labor productivity, taken weekly, for each fulfillment center (FC) between 2011 and 2020. We also have monthly data on turnover for each FC between 2017 and 2020, which we will lag by 1 month because we want to observe how turnover from the previous month affects current productivity. We aggregated the weekly labor productivity data into monthly data by grouping by FC and month and taking the average values per FC per month, merging it with the monthly turnover data to get one dataset with the following features:

|  | Data Field | Data Description |
|---|---|---|
| 1 | fc_id | unique identifier for each fulfillment center |
| 2 | year_month | month of current observation in YYYY-MM format |
| 3 | turn_full_last_month | turnover of full-time employees from previous month |
| 4 | turn_part_last_month | turnover of part-time employees from previous month |
| 5 | turn_sum_last_month | turnover of full and part-time employees from previous month |
| 6 | boxes_per_hr | boxes moved per hour, the primary measure of productivity |

| 7 | bph_actual_minus_fcast | difference between actual and forecasted boxes per hour |
|---|---|---|
| 8 | cargo_loss | boxes lost |
| 9 | damages_per_sales | $ value of damages incurred as a percentage of total $ value of sales |
| 10 | fcast_cargo_loss | forecasted cargo loss |
| 11 | fcast_damages_per_sales | forecasted damages value as a percentage of sales |
| 12 | fcast_labor_per_sales | forecasted labor value as a percentage of sales |
| 13 | fcast_ot_hrs_pct_tot_prod_hrs | forecasted overtime hours as a percentage of total hours |
| 14 | fcast_prod_avg_wage | forecasted average wage |
| 15 | fcast_receiving_cartons | forecasted cartons received |
| 16 | fcast_vol_boxes_wkly | forecasted boxes processed per week |
| 17 | ft_headcount_wkly | forecasted headcount per week |
| 18 | labor_per_box | $ value of labor per box handled (total wage bill / total boxes) |
| 19 | labor_per_sales | $ value of labor as a percentage of sales (total wage bill / total sales) |
| 20 | leave_behinds | boxes that were left behind and didn't make it onto a truck |
| 21 | missing_wrong | boxes that were missing or wrong |
| 22 | ot_hrs | total overtime hours worked |
| 23 | ot_hrs_pct_tot_prod_hrs | overtime hours as a percentage of total hours worked |
| 24 | prod_avg_wage | average wage of workers |
| 25 | pt_headcount_wkly | headcount of part-time workers per week |
| 26 | pt_ot_hrs | overtime hours worked by part-time workers |
| 27 | pt_tot_hrs | total hours worked by part-time workers |
| 28 | receiving_cartons | cartons received per week |
| 29 | temp_headcount_wkly | headcount of temp workers per week |
| 30 | temp_hrs | total hours worked by temp workers |
| 31 | temp_ot_hrs | overtime hours worked by temp (3rd party contracted) |

| 32 | tot_damages | $ value of total damages incurred |
|---|---|---|
| 33 | tot_headcount_wkly | total headcount of all workers |
| 34 | tot_prod_hrs | total hours worked by all workers |
| 35 | tot_sales | total sales |
| 36 | vol_boxes_wkly | boxes processed per week |

In the labor productivity data, there exists a measure of forecasted productivity (this is the target level that the FCs aim to reach), which can be used as a measure of how close the FCs are to their optimal productivity.

## 4. Data Exploration

### 4.1 Data Preprocessing

We started with 2 separate data sets. The first dataset is the turnover dataset and second dataset is the labor productivity dataset. They contain the following:
- **Turnover:** turnover numbers—how many people left their job, either voluntarily or involuntarily (laid off or fired)—at each FC (fulfillment center) from 2017 to 2020 for both part time and full time workers. This data is in the wide format (rows: FC number; columns: month)
- **Labor productivity (l/p):** l/p values for each week for each FC from 2011-2021. This data is in the long format (rows: FC and month; columns: measures of various l/p values such as headcount, cargo loss, productivity, etc)

Ultimately, we converted the turnover data to the long format so that we could work with both sets with the same procedure. Then, we lagged the turnover data so that for each observation, we know the turnover values from the previous month. We do this because for each observation, we know the turnover values from the previous month. We want to use the last month's turnover values to predict this month's productivity so that knowing last month's turnover, we can make decisions for this month to maximize productivity.

We then added a new column that converted the labor productivity week column (we_date, which indicates the date of the end of the week in YYYY-MM-DD format) to YYYY-MM format. We converted all dates in both datasets to the same format. Specifically,

datetime objects. For the l/p data, these dates came in unix format -- integer number of days since unix epoch, whereas for turnover data, the dates were originally strings.

We grouped the l/p data by the FC ID, which identifies each specific fulfillment center, as well as month, and took the average value of each of these groups to convert the weekly data to monthly data. The values represent the average value for that FC for each respective month. For example, all l/p observations are the monthly averages for each FC, or the average across all weekly observations for that month for that respective FC.

Then, we merged the turnover with labor productivity data on FC ID and month. The exact intuition for how we did this is as follows. There is much more l/p data, which spans from 2011-2021, than there is turnover data, which spans from 2017-2020. Ultimately, we only merged with inner join because we only want to keep observations for which we have both turnover and l/p data present. We do not care about the observations for which we don't have any turnover data -- it is not reasonable to interpolate these missing values from other observations because it's very hard to say what factors turnover depends on, and every year could be different from another. Furthermore, we don't have enough years of turnover data to look at this time series analysis. Since turnover is the explanatory variable, we do not want our own biases in an interpolation to impact the model results.
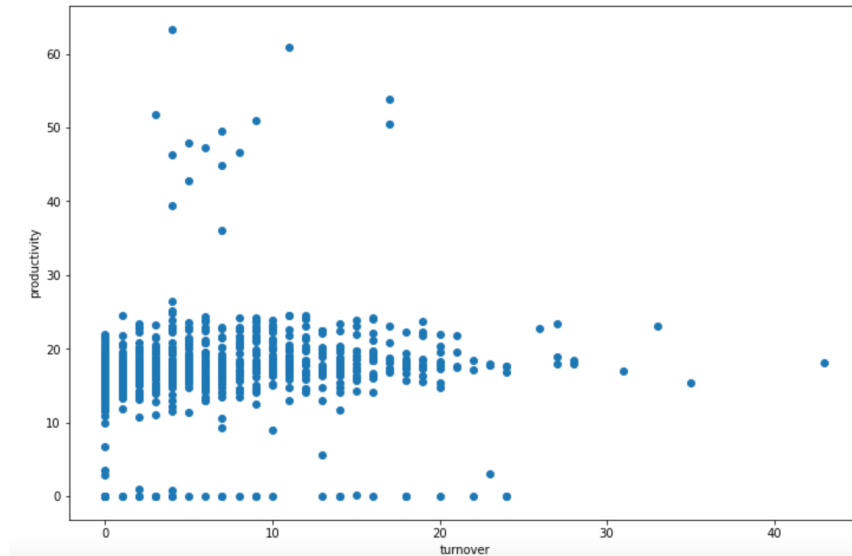
Additionally, we also do not care about observations for which we don't have l/p data—there is too much missing data here. After completing the inner join, there are a little over 1000 observations left.

## 4.2 Exploratory Data Analysis

After we completed our data preprocessing, we plotted turnover (FT (full time) and PT (part time)) versus productivity (boxes_per_hr). Boxes_per_hr is the primary measurement of productivity within FCs.

The results were not very revealing; there is no clear correlation between turnover and productivity, with the majority of points being clustered between values of 0 and 20 for turnover and 10 and 30 for productivity.

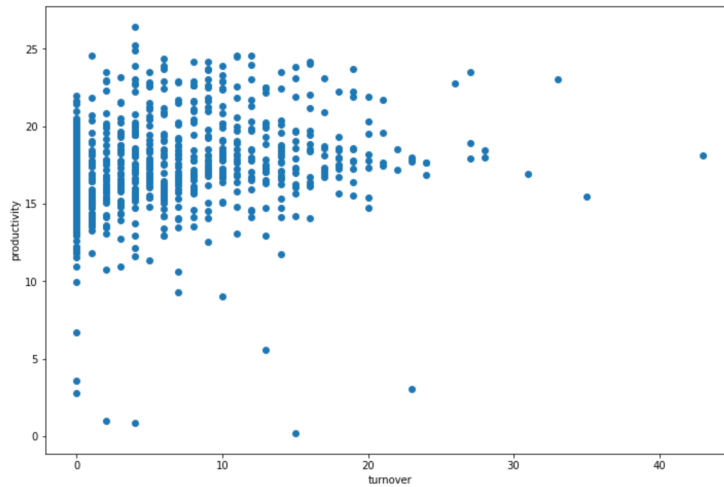*Simple regression of productivity on FT and PT turnover*

| Dep. Variable: | boxes_per_hr | R-squared (uncentered): | 0.394 |
|---|---|---|---|
| Model: | OLS | Adj. R-squared (uncentered): | 0.392 |
| Method: | Least Squares | F-statistic: | 338.5 |
| Date: | Mon, 29 Nov 2021 | Prob (F-statistic): | 5.00e-114 |
| Time: | 18:38:50 | Log-Likelihood: | -4245.9 |
| No. Observations: | 1045 | AIC: | 8496. |
| Df Residuals: | 1043 | BIC: | 8506. |
| Df Model: | 2 | | |
| Covariance Type: | nonrobust | | |

| | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| turn_full_last_month | 1.5007 | 0.108 | 13.929 | 0.000 | 1.289 | 1.712 |
| turn_part_last_month | 1.0848 | 0.246 | 4.410 | 0.000 | 0.602 | 1.567 |

| | | | |
|---|---|---|---|
| Omnibus: | 208.199 | Durbin-Watson: | 0.903 |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 625.519 |
| Skew: | -0.987 | Prob(JB): | 1.48e-136 |
| Kurtosis: | 6.235 | Cond. No. | 3.79 |

We have a low $R^2$ value here, 0.394. Ultimately, many outliers are probably contributing to low the $R^2$ value that we have here. First of all, we try removing outliers, observations where productivity is greater than 30 and productivity is equal to 0, and try the regression again:

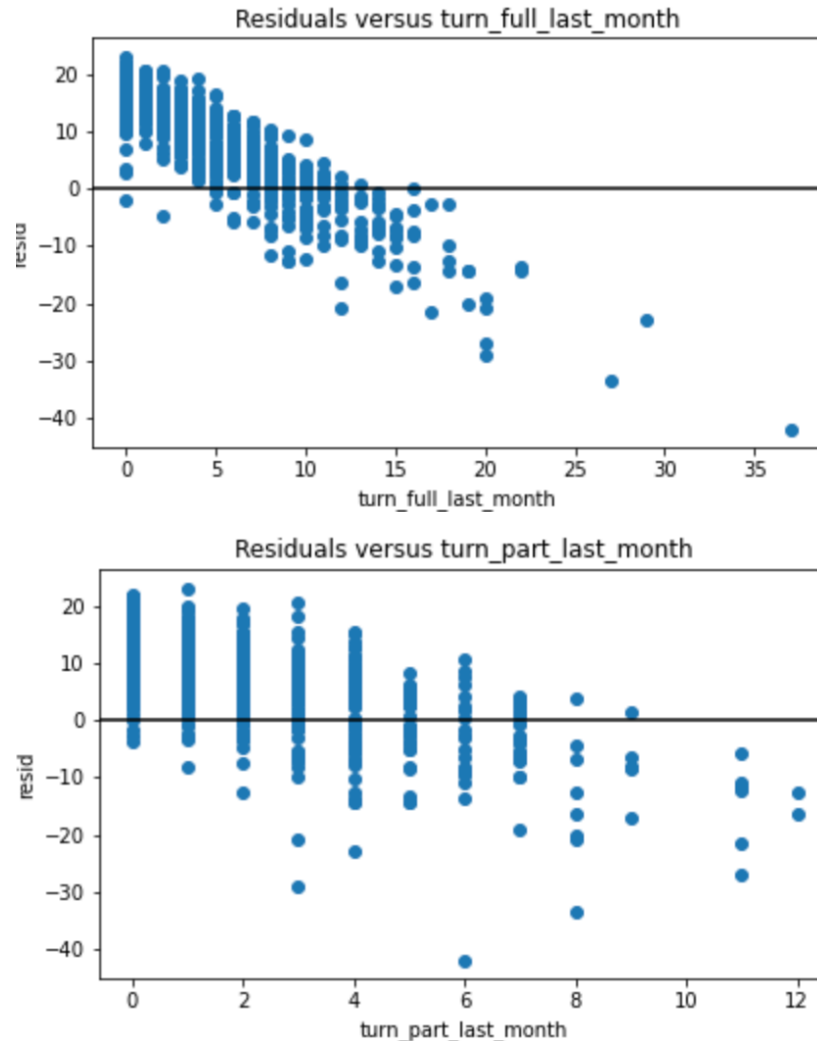*Simple regression of productivity on FT and PT turnover (removing outliers)*



| Dep. Variable: | boxes_per_hr | R-squared (uncentered): | 0.421 |
|---|---|---|---|
| Model: | OLS | Adj. R-squared (uncentered): | 0.420 |
| Method: | Least Squares | F-statistic: | 355.1 |
| Date: | Mon, 29 Nov 2021 | Prob (F-statistic): | 1.27e-116 |
| Time: | 18:26:13 | Log-Likelihood: | -3931.8 |
| No. Observations: | 979 | AIC: | 7868. |
| Df Residuals: | 977 | BIC: | 7877. |
| Df Model: | 2 | | |
| Covariance Type: | nonrobust | | |

| | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| turn_full_last_month | 1.3851 | 0.108 | 12.781 | 0.000 | 1.172 | 1.598 |
| turn_part_last_month | 1.4632 | 0.245 | 5.970 | 0.000 | 0.982 | 1.944 |

| Omnibus: | 283.613 | Durbin-Watson: | 0.722 |
|---|---|---|---|
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 729.094 |
| Skew: | -1.512 | Prob(JB): | 4.78e-159 |
| Kurtosis: | 5.955 | Cond. No. | 3.81 |

The $R^2$ is still low because we need to include other features. Below, we plot the residual plots:

*Residual plots*

Residuals versus turn_full_last_month



Residuals versus turn_part_last_month

Examining the residual plots, we see that there is clearly a non-linear relationship with FT turnover and possibly PT turnover as well because the residuals are not randomly distributed. To resolve this issue of non-linearity, we introduce $turnover^2$. In total, we include 2 new features -- $turnover_{FT}^2$ and $turnover_{PT}^2$, for FT and PT, respectively.

*Simple regression of productivity on FT and PT turnover (removing outliers, adding in $turnover^2$)*

| | | | |
|---|---|---|---|
| **Dep. Variable:** | boxes_per_hr | **R-squared (uncentered):** | 0.520 |
| **Model:** | OLS | **Adj. R-squared (uncentered):** | 0.518 |
| **Method:** | Least Squares | **F-statistic:** | 262.1 |
| **Date:** | Mon, 29 Nov 2021 | **Prob (F-statistic):** | 1.64e-152 |
| **Time:** | 18:58:23 | **Log-Likelihood:** | -3810.6 |
| **No. Observations:** | 971 | **AIC:** | 7629. |
| **Df Residuals:** | 967 | **BIC:** | 7649. |
| **Df Model:** | 4 | | |
| **Covariance Type:** | nonrobust | | |

| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| **turn_full_last_month** | 2.6654 | 0.193 | 13.797 | 0.000 | 2.286 | 3.044 |
| **turn_part_last_month** | 3.4589 | 0.517 | 6.691 | 0.000 | 2.445 | 4.473 |
| **turn_full_sq** | -0.0966 | 0.010 | -10.119 | 0.000 | -0.115 | -0.078 |
| **turn_part_sq** | -0.3547 | 0.059 | -6.003 | 0.000 | -0.471 | -0.239 |

| | | | |
|---|---|---|---|
| **Omnibus:** | 86.579 | **Durbin-Watson:** | 0.903 |
| **Prob(Omnibus):** | 0.000 | **Jarque-Bera (JB):** | 51.652 |
| **Skew:** | -0.428 | **Prob(JB):** | 6.08e-12 |
| **Kurtosis:** | 2.262 | **Cond. No.** | 126. |

We conclude that the residual plots look slightly better, but still not great. This can ultimately probably get resolved by adding in more features into the regression.

## 4.3 Model 1: Linear regression with entire dataset

For our response variable, we used productivity for this month (boxes_per_hr). Our explanatory variables were turnover from last month—full time turnover, part time turnover, $turnover_{FT}^2$ and $turnover_{PT}^2$. We used a correlation matrix to determine which additional features to include in the regression. We left out variables that are highly correlated with each other to avoid multicollinearity and to reduce the complexity of the model.
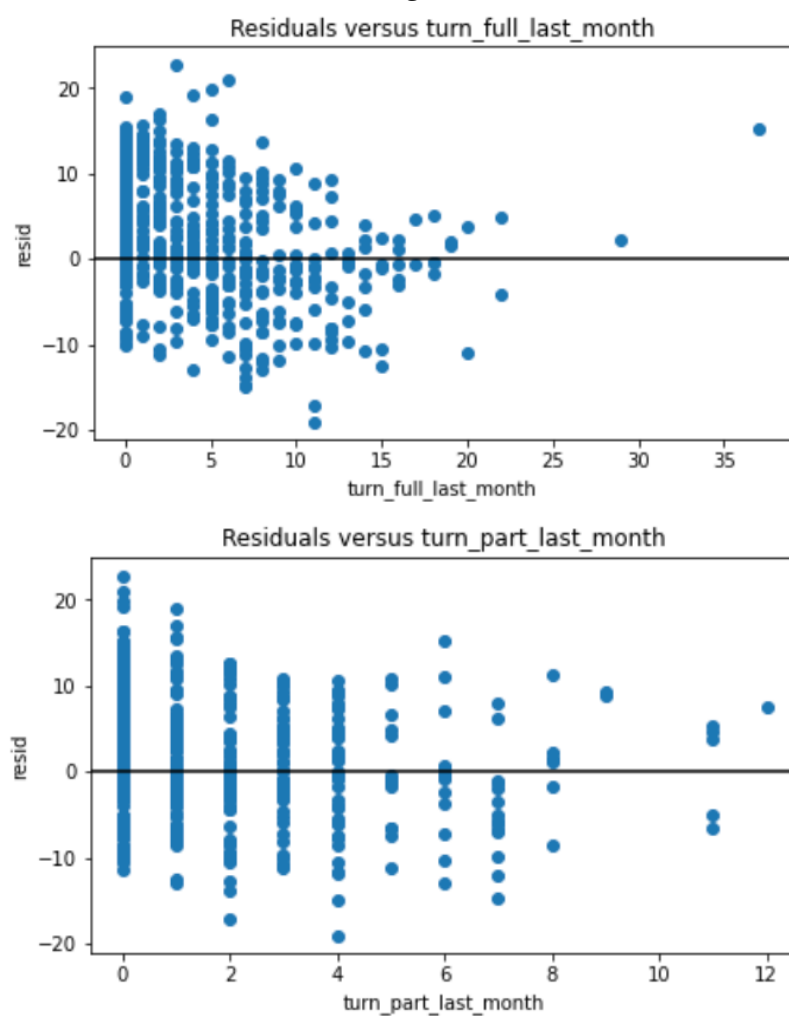
| | fc_id | turn_full_last_month | turn_part_last_month | turn_sum_last_month | boxes_per_hr | bph_actual_minus_fcast | cargo_loss | dam |
|---|---|---|---|---|---|---|---|---|
| fc_id | 1.000000 | -0.134062 | -0.090215 | -0.131370 | 0.091614 | 0.005249 | -0.107657 | |
| turn_full_last_month | -0.134062 | 1.000000 | 0.621050 | 0.961800 | 0.184257 | 0.043074 | -0.018200 | |
| turn_part_last_month | -0.090215 | 0.621050 | 1.000000 | 0.811887 | 0.190507 | -0.054227 | -0.021081 | |
| turn_sum_last_month | -0.131370 | 0.961800 | 0.811887 | 1.000000 | 0.203789 | 0.013104 | -0.020920 | |
| boxes_per_hr | 0.091614 | 0.184257 | 0.190507 | 0.203789 | 1.000000 | 0.203104 | -0.039547 | |
| bph_actual_minus_fcast | 0.005249 | 0.043074 | -0.054227 | 0.013104 | 0.203104 | 1.000000 | -0.005394 | |
| cargo_loss | -0.107657 | -0.018200 | -0.021081 | -0.020920 | -0.039547 | -0.005394 | 1.000000 | |
| damages_per_sales | 0.021406 | 0.071186 | 0.022703 | 0.060955 | -0.220727 | -0.012465 | -0.001251 | |
| fcast_cargo_loss | -0.007631 | -0.032107 | 0.040849 | -0.008623 | -0.164649 | -0.413897 | -0.034151 | |
| fcast_damages_per_sales | -0.052667 | -0.028270 | -0.017384 | -0.027096 | -0.076604 | -0.029526 | 0.000982 | |
| fcast_labor_per_sales | -0.087851 | 0.074474 | 0.013778 | 0.060286 | -0.369363 | 0.582407 | 0.033059 | |
| fcast_ot_hrs_pct_tot_prod_hrs | 0.113858 | -0.052233 | 0.011509 | -0.034873 | -0.296021 | -0.328835 | -0.021582 | |
| fcast_prod_avg_wage | -0.154537 | 0.022464 | 0.071550 | 0.041767 | -0.031524 | -0.589841 | -0.000033 | |
| fcast_receiving_cartons | -0.044098 | -0.026786 | -0.037488 | -0.033067 | 0.165133 | -0.091910 | -0.047280 | |
| fcast_vol_boxes_wkly | -0.075993 | -0.028637 | -0.041537 | -0.035862 | 0.194837 | -0.101613 | -0.048496 | |
| ft_headcount_wkly | -0.073175 | -0.093518 | -0.093301 | -0.102248 | 0.078572 | -0.116306 | -0.049018 | |
| labor_per_box | 0.028176 | -0.007872 | 0.008340 | -0.002951 | -0.244143 | 0.007983 | -0.000216 | |
| labor_per_sales | 0.035896 | 0.045831 | 0.024627 | 0.042740 | -0.267028 | -0.019487 | -0.001263 | |
| leave_behinds | -0.008878 | -0.023641 | 0.000958 | -0.017275 | -0.246778 | -0.212251 | -0.009640 | |
| missing_wrong | 0.025856 | -0.010363 | -0.020325 | -0.014818 | -0.165989 | -0.005439 | -0.001820 | |

-

The additional features that we selected are:
- **bph_actual_minus_fcast** (difference between the actual and forecasted boxes per hour)
- **damages_per_sales** ($ value of damages incurred as a percentage of total $ value of sales)
- **ft_headcount_wkly** (full time headcount -- because we aggregated weekly to monthly data, these values are the average FT headcount across all the weeks in each month)
- **leave_behinds** (boxes that were left behind and didn't make it onto a truck)
- **missing_wrong** (boxes that were missing or defective in some way)
- **ot_hrs** (total overtime hours worked in that FC)
- **pt_headcount_wkly** (headcount of part-time workers per week)
- **temp_headcount_wkly** (headcount of temporary (e.g. contracted) workers per week)
- **tot_sales** (total sales)

Once selecting the variables, we split the data into training and testing sets, with 70% of the data in training and 30% of the data in testing. We created an OLS model built on the training data:

*OLS Model on Training Data*



Residuals versus turn_full_last_month



Residuals versus turn_part_last_month

OLS Regression Results

| Dep. Variable: | boxes_per_hr | R-squared (uncentered): | 0.807 |
|---|---|---|---|
| Model: | OLS | Adj. R-squared (uncentered): | 0.803 |
| Method: | Least Squares | F-statistic: | 214.0 |
| Date: | Mon, 29 Nov 2021 | Prob (F-statistic): | 1.32e-227 |
| Time: | 19:01:54 | Log-Likelihood: | -2352.6 |
| No. Observations: | 679 | AIC: | 4731. |
| Df Residuals: | 666 | BIC: | 4790. |
| Df Model: | 13 | | |
| Covariance Type: | nonrobust | | |

| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| turn_full_last_month | 1.1632 | 0.151 | 7.680 | 0.000 | 0.866 | 1.461 |
| turn_part_last_month | 1.5696 | 0.405 | 3.871 | 0.000 | 0.773 | 2.366 |
| turn_full_sq | -0.0373 | 0.007 | -5.298 | 0.000 | -0.051 | -0.024 |
| turn_part_sq | -0.1451 | 0.046 | -3.133 | 0.002 | -0.236 | -0.054 |
| bph_actual_minus_fcast | -0.3045 | 0.089 | -3.439 | 0.001 | -0.478 | -0.131 |
| damages_per_sales | -5.3276 | 5.337 | -0.998 | 0.318 | -15.806 | 5.151 |
| ft_headcount_wkly | 0.0838 | 0.012 | 6.707 | 0.000 | 0.059 | 0.108 |
| leave_behinds | 26.7445 | 133.787 | 0.200 | 0.842 | -235.950 | 289.439 |
| missing_wrong | 13.9003 | 10.359 | 1.342 | 0.180 | -6.440 | 34.240 |
| ot_hrs | -0.0026 | 0.001 | -1.978 | 0.048 | -0.005 | -1.94e-05 |
| pt_headcount_wkly | 0.0249 | 0.030 | 0.841 | 0.401 | -0.033 | 0.083 |
| temp_headcount_wkly | 0.0223 | 0.021 | 1.051 | 0.294 | -0.019 | 0.064 |
| tot_sales | -2.49e-07 | 3.77e-07 | -0.660 | 0.509 | -9.9e-07 | 4.91e-07 |

| Omnibus: | 12.279 | Durbin-Watson: | 1.685 |
|---|---|---|---|
| Prob(Omnibus): | 0.002 | Jarque-Bera (JB): | 10.068 |
| Skew: | -0.218 | Prob(JB): | 0.00651 |
| Kurtosis: | 2.592 | Cond. No. | 2.77e+09 |

We conclude that the features with the highest significance are the turnover variables, bph_actual_minus_fcast, ft_headcount_wkly, and ot_hrs. This makes sense because the difference between actual and predicted productivity will obviously be highly correlated with actual productivity. Headcount also impacts how many boxes per hr an FC can process; if many people are working overtime, they may be more exhausted and less productive, which is indicated by the negative relationship between ot_hrs and boxes_pr_hr.

The residuals are much more randomly distributed in this model. We observe a RMSE of 7.978 for the testing data and a RMSE of 7.736 for the training data. These are similar values, so it is unlikely to be overfitting the data.
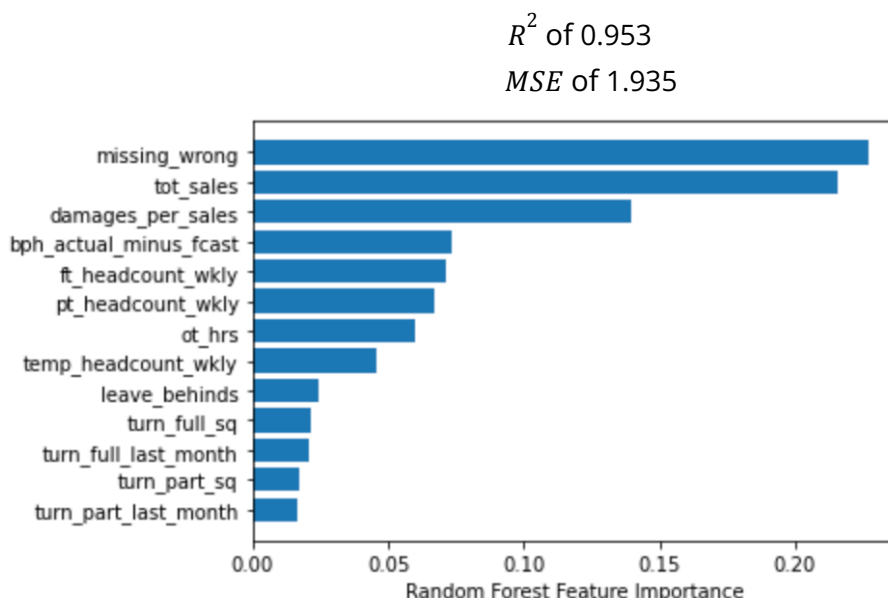
*OLS results on entire data set (no train/test split):*

| | | | |
|---|---|---|---|
| **Dep. Variable:** | boxes_per_hr | **R-squared (uncentered):** | 0.807 |
| **Model:** | OLS | **Adj. R-squared (uncentered):** | 0.804 |
| **Method:** | Least Squares | **F-statistic:** | 308.4 |
| **Date:** | Mon, 29 Nov 2021 | **Prob (F-statistic):** | 0.00 |
| **Time:** | 19:04:59 | **Log-Likelihood:** | -3368.2 |
| **No. Observations:** | 971 | **AIC:** | 6762. |
| **Df Residuals:** | 958 | **BIC:** | 6826. |
| **Df Model:** | 13 | | |
| **Covariance Type:** | nonrobust | | |

| | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| turn_full_last_month | 1.2356 | 0.131 | 9.464 | 0.000 | 0.979 | 1.492 |
| turn_part_last_month | 1.7625 | 0.336 | 5.240 | 0.000 | 1.102 | 2.423 |
| turn_full_sq | -0.0422 | 0.006 | -6.677 | 0.000 | -0.055 | -0.030 |
| turn_part_sq | -0.1661 | 0.038 | -4.344 | 0.000 | -0.241 | -0.091 |
| bph_actual_minus_fcast | -0.1242 | 0.067 | -1.852 | 0.064 | -0.256 | 0.007 |
| damages_per_sales | -4.2665 | 5.321 | -0.802 | 0.423 | -14.708 | 6.175 |
| ft_headcount_wkly | 0.0878 | 0.010 | 8.381 | 0.000 | 0.067 | 0.108 |
| leave_behinds | -21.2036 | 110.518 | -0.192 | 0.848 | -238.089 | 195.682 |
| missing_wrong | 17.7212 | 10.361 | 1.710 | 0.088 | -2.612 | 38.054 |
| ot_hrs | -0.0026 | 0.001 | -2.471 | 0.014 | -0.005 | -0.001 |
| pt_headcount_wkly | 0.0347 | 0.025 | 1.385 | 0.167 | -0.014 | 0.084 |
| temp_headcount_wkly | 0.0408 | 0.017 | 2.360 | 0.018 | 0.007 | 0.075 |
| tot_sales | -4.502e-07 | 3.15e-07 | -1.430 | 0.153 | -1.07e-06 | 1.68e-07 |

| | | | |
|---|---|---|---|
| **Omnibus:** | 15.581 | **Durbin-Watson:** | 1.046 |
| **Prob(Omnibus):** | 0.000 | **Jarque-Bera (JB):** | 13.753 |
| **Skew:** | -0.233 | **Prob(JB):** | 0.00103 |
| **Kurtosis:** | 2.648 | **Cond. No.** | 2.75e+09 |

As deduced earlier, we have a RMSE of 7.810 -- which yields similar results as the train-test split.

## 4.4 Model 2: Random Forest Regression using 30% train/test split

The next model we try is a random forest regression using 30% train/test split. Bph_actual_minus_fcast had 6 nan values and the model couldn't parse these values, so had to re-clean data and remove rows with nan in them. After doing so, we get the following results:

$$R^2 \text{ of } 0.953$$
$$MSE \text{ of } 1.935$$



Interestingly, feature importances actually rank turnover variables as least important in the random forest approach.

## 4.5 OLS vs RFR

When comparing the two previous methods we used, the OLS has more interpretable results, but it is not as accurate. From the OLS coefficients, we are able to see the marginal effects of each feature on productivity. Meanwhile, the RFR is much more accurate, but it doesn't have nicely interpretable results like OLS. It also doesn't know how to extrapolate outside of the data set, which could be problematic if, for example, we have new FC turnover levels that are outside of the range of turnover that was present in the training data (or in any case where the new data has observations that fall out of the range used in training). The RFR probably performs better since there isn't a linear relationship between turnover and productivity and it is able to identify nonlinear relationships, but there is a risk that it could be overfitting.

If we wanted a more accurate model that is able to extrapolate, we could train a neural network. But this would ideally be done with more data. Unfortunately, 1000 observations is not enough to do this.

## 4.6 Model 3: Separate model for each FC

Finally, the third model we try is an optimization where we make a separate model for each FC. Unfortunately, each FC only has around 40 observations, which is not enough data to make a good model. By including features like headcount in the model that uses the entire dataset, we are able to account for differences in FC size.

**Optimization**

We aimed to maximize the productivity for a given FC for a given month. We used coefficients from the train-test-split model.

***Objective Function:***

*maximize productivity for a given FC for a given month*

$max\ 1.163F\ +\ 1.570P\ -\ 0.037F^2\ -\ 0.145P^2$

*F* and *P* are the levels of full time and part time turnover per month (e.g. how many people should leave the FC to maximize turnover). A negative turnover implies that people should be hired.

***Constraints (specific to each FC for each month):***

Turnover for next month cannot exceed the headcount that the FC starts out with (equal to headcount for the last week of the previous month). We used the data from the previous month to set this constraint.

*Example with FC 6710:*

The most recent data we have is for 2020-05. We can find levels of turnover that will maximize productivity for 2020-06. In the last week of 2020-05, full-time headcount was 187.0 and part-time headcount was 25.0. Based on the below results, the optimal level of turnover is 16 FT and 5 PT for a max level of productivity of 13.361 boxes_per_hr (not accounting for other features).

This objective function does not take into account the constants from the other features such as ot_hrs and missing_wrong. These features could be included in the function by multiplying the average value of that feature from 2020-05 by the coefficient

determined by the model, but this would just serve as a constant that shifts the objective function—the values of F and P that maximize the function would remain the same.

```
In [3]:    1  m = Model(Gurobi.Optimizer);

           Academic license - for non-commercial use only - expires 2022-09-06
```

```
In [4]:    1  # variable for Full time turnover
           2  @variable(m, F <= 187, Int)
           3  # variable for Part time turnover
           4  @variable(m, P <= 25, Int)
           5  @objective(m, Max, 1.163*F + 1.570*P - 0.037*F^2 - 0.145*P^2)
```

Out[4]: $F$

```
In [10]:   1  print(m)
```

$$\begin{aligned} \max \quad & -0.037F^2 - 0.145P^2 + 1.163F + 1.57P \\ \text{Subject to} \quad & F \leq 187.0 \\ & P \leq 25.0 \\ & F \in \mathbb{Z} \\ & P \in \mathbb{Z} \end{aligned}$$

```
In [7]:    1  optimize!(m)
```

...

```
In [9]:    1  @show value(F)
           2  @show value(P)
           3
           4  @show objective_value(m);

           value(F) = 16.0
           value(P) = 5.0
           objective_value(m) = 13.361000000000004
```

## 5. Results & Insight

  Having some amount of turnover actually improves productivity—on one hand, this may feel harsh at first because it seems like the warehouses perform better when workers leave. On the other hand, it may be the case that many warehouses have too many workers, leading to inefficiencies ("too many cooks in the kitchen"), and some may actually be better off if they leave.

  The social science aspect of this study would look at which "workers are the ones that should go?" "Is it ethical to lay off workers to increase productivity?" "Are there other issues going on in the warehouses that are causing the inefficiencies that arise with having

lots of workers -- e.g. is there a way to increase productivity without laying people off?" Turnover was just one of the features (and as the Random Forest Regression showed, it was not the most important one) that affects productivity, so maybe there are other areas of focus that the warehouses should be targeting first to increase productivity.

However, turnover is usually associated with negative worker well-being. It doesn't feel right to say that turnover is a good thing because usually workers have a relatively inelastic supply of labor -- they would rather stay at their current job instead of going through the trouble of finding a new one, unless the conditions in their current job are bad.

Ultimately, it would be great to have more turnover data for a longer time period to build a stronger model. There may even be a seasonality component within turnover/productivity that we are unable to see with the current small dataset. It will be interesting to take these results back to the research group and see what they believe is the next best area to focus on. The group's goal is to improve warehouse worker wellbeing, not necessarily warehouse productivity. But higher levels of productivity may be correlated with higher levels of wellbeing. However, this is tricky because generally higher turnover is correlated with lower levels of wellbeing. Is there a way to improve productivity without negatively impacting well being through higher turnover?

Potentially, though, turnover may not actually be bad for wellbeing. For example, there may be people who would have left the job anyways because they have found a better job or they are going through another transition in life. These people probably should not be lumped together with the people who are leaving because they are unhappy with their current job. Perhaps there are enough of these people leaving each month to increase productivity -- in this case, turnover would not be a bad thing because they are not leaving for reasons related to negative well-being. Clearly, this is a complicated problem to approach both analytically and ethically. However, the results we obtained will certainly provide valuable insight to the research group.